



Available at
www.ElsevierComputerScience.com
POWERED BY SCIENCE @ DIRECT®

Journal of Computer and System Sciences 69 (2004) 196–234

JOURNAL of
COMPUTER
AND SYSTEM
SCIENCES

<http://www.elsevier.com/locate/jcss>

On approximating weighted sums with exponentially many terms[☆]

Deepak Chawla,¹ Lin Li, and Stephen Scott*

Department of Computer Science, University of Nebraska, Lincoln, NE 68588-0115, USA

Received 28 March 2003; revised 8 January 2004

Abstract

Multiplicative weight-update algorithms such as Winnow and Weighted Majority have been studied extensively due to their on-line mistake bounds' logarithmic dependence on N , the total number of inputs, which allows them to be applied to problems where N is exponential. However, a large N requires techniques to efficiently compute the weighted sums of inputs to these algorithms. In special cases, the weighted sum can be exactly computed efficiently, but for numerous problems such an approach seems infeasible. Thus we explore applications of Markov chain Monte Carlo (MCMC) methods to estimate the total weight. Our methods are very general and applicable to any representation of a learning problem for which the inputs to a linear learning algorithm can be represented as states in a completely connected, untruncated Markov chain. We give theoretical worst-case guarantees on our technique and then apply it to two problems: learning DNF formulas using Winnow, and pruning classifier ensembles using Weighted Majority. We then present empirical results on simulated data indicating that in practice, the time complexity is much better than what is implied by our worst-case theoretical analysis.

© 2003 Elsevier Inc. All rights reserved.

Keywords: Markov chain Monte Carlo approximation; Winnow; Weighted Majority; Multiplicative weight updates; Perceptron; DNF learning; Boosting

1. Introduction

Multiplicative weight-update algorithms (e.g. [6,21,24]) have been studied extensively due to their on-line mistake bounds' logarithmic dependence on N , the total number of inputs. (These

[☆] A preliminary version [7] of this paper appeared in COLT 2001.

*Corresponding author.

E-mail address: sscott@cse.unl.edu (S. Scott).

URL: <http://www.cse.unl.edu/~sscott>.

¹ Now at EMC Corporation, Raleigh-Durham, NC.

bounds can be translated into PAC sample complexity bounds via a simple procedure [22].) This *attribute efficiency* allows them to be applied to problems where N is exponential in the input size, which is the case in many applications, including using Winnow [21] to learn DNF formulas in unrestricted domains and using the Weighted Majority algorithm (WM [24]) to predict nearly as well as the best pruning of a classifier ensemble (from e.g. boosting). However, a large N requires techniques to efficiently compute the weighted sums of inputs to WM and Winnow. One method of doing this is to exploit commonalities among the inputs, partitioning them into a polynomial number of *groups* such that given a single member of each group, the total weight contribution of that group can be efficiently computed [13–15,25,30,39]. But many WM and Winnow applications do not appear to exhibit such structure, so it seems that a brute-force implementation is the only option to guarantee complete correctness.² Thus we explore applications of Markov chain Monte Carlo (MCMC) methods to estimate the total weight without the need for special structure in the problem. Our methods are very general and applicable to any representation of a learning problem for which the inputs to the linear learning algorithm can be represented as states in a completely connected, untruncated³ Markov chain. In this paper we apply our results to two such problems, described below.

First we study learning DNF formulas (e.g. [5]) using Winnow⁴ [21] and not using membership queries. We enumerate all possible DNF terms and use Winnow to learn a monotone disjunction over these terms, which it can do while making $O(K \log N)$ prediction mistakes, where K is the number of relevant terms and N is the total number of terms. So a brute-force implementation of Winnow makes a polynomial number of errors on arbitrary examples (i.e. with no distributional assumptions) and does not require membership queries. However, a brute-force implementation requires exponential time to compute the weighted sum of the inputs. So we apply our MCMC-based results to estimate this sum.

Next we investigate pruning a classifier ensemble (from e.g. boosting), which can reduce overfitting and time for evaluation [26,40]. We use the Weighted Majority algorithm (WM) [24], using all possible prunings as experts. WM is guaranteed to not make many more prediction mistakes than the best expert, so we know that a brute-force WM will perform nearly as well as the best pruning. However, the exponential number of prunings motivates us to use an MCMC approach to approximate the weighted sum of the experts' predictions.

MCMC methods [16] have been applied to problems in approximate summation, where the goal is to approximate $W = \sum_{x \in \Omega} s(x)$, where s is a positive function and Ω is a finite set of combinatorial structures. It involves defining an ergodic Markov chain \mathcal{M} with state space Ω and stationary distribution π . Then one repeatedly simulates \mathcal{M} to draw samples almost according to π . Under appropriate conditions, this technique yields accuracy guarantees. E.g. sometimes one can guarantee that the estimate of the sum is within a factor ε of the true value (with high

²For additive weight-update algorithms such as the Perceptron algorithm, kernels can often be used to exactly compute the weighted sums (e.g. [8,29,36]), though a kernel function might not exist for the desired mapping of features.

³We believe that the untruncated requirement can be removed by generalizing results of Morris and Sinclair [28] (see Section 4.2).

⁴We also study the application of our methods to learning DNF via Rosenblatt's Perceptron [33] algorithm, though this is done only for contrast with Winnow since exact sums for DNF learning via Perceptron can be computed with kernels [18].

probability). When this is true and the estimation algorithm requires only polynomial time, the algorithm is called a *fully polynomial randomized approximation scheme* (FPRAS).

We combine two FPRASs for application to estimating weighted sums. The first approximator is for the *approximate knapsack problem* [10,28], where given a positive real vector \vec{w} and real number b , the goal is to estimate $|\{\vec{p} \in \{0,1\}^n: \vec{w} \cdot p(\vec{x}) \leq b\}|$ within a multiplicative factor ε . The other FPRAS is for estimating the sum of the weights of a *weighted matching* of a graph: for a graph G and $\lambda \geq 0$, approximate $Z_G(\lambda) = \sum_{k=0}^n m_k \lambda^k$, where m_k is the number of matchings in G of size k and n is the number of nodes. This problem has applications to the *monomer-dimer problem* of statistical physics [16].

While we have thoroughly analyzed our approach in the context of these two problems, our results do not guarantee efficient algorithms for learning DNF and for finding the best pruning.⁵ But we do provide theoretical machinery that could potentially be applied to analyze algorithms that learn e.g. restricted cases of DNF, including subclasses of DNF formulas and/or specific distributions over examples. Further, our experimental results provide interesting insights into the algorithms' behaviors and show that the weighted sums can be approximated well despite the pessimistic worst-case bounds. Couple this with the fact that good approximations of the weighted sums are not always necessary to accurately simulate Winnow and WM (since we are only interested in the predictions made based on these weighted sums, not the sums themselves), and our results have potential to be effective tools in theory and in practice.

The rest of this paper is organized as follows. In Section 2 we give background on the on-line learning model and summarize related work in learning DNF formulas, pruning ensembles of classifiers, and MCMC methods. Section 3 presents our algorithm and Markov chain, and proves general bounds on the accuracy and time complexity of our estimation procedure. In Section 4 we apply these results to the problems of learning DNF formulas with Winnow and Perceptron and pruning ensembles with Weighted Majority. Then some empirical results appear in Section 5. Finally, we conclude in Section 6 with a description of future and ongoing work.

2. Related work

2.1. The on-line learning model

We focus on *on-line* learning algorithms, where learning proceeds in a series of *trials*.⁶ In trial t , an example \vec{X}_t is presented to the learning algorithm A , which makes a prediction⁷ $\hat{\ell}_t$ of \vec{X}_t 's label. After this prediction is made, A is told the true label ℓ_t , which A uses to update its hypothesis before making future predictions. If $\ell_t \neq \hat{\ell}_t$, we say that A made a *prediction mistake*. If M is the set of examples for which a mistake is made, the goal is to minimize $|M|$ on any sequence of adversarially-generated examples $\mathcal{X} = (\vec{X}_1, \dots, \vec{X}_\tau)$. Below we overview the on-line learning algorithms Winnow [21], Perceptron [33], and Weighted Majority (WM) [24].

⁵This is not surprising since it is unlikely that an efficient distribution-free DNF-learning algorithm exists [2,3].

⁶In Sections 3 and 4, our results focus on only the current trial, so we omit the subscript t unless it is not clear from context.

⁷We assume $\ell_t, \hat{\ell}_t \in \{-1, +1\}$.

Winnnow maintains a weight vector $\vec{w} \in \mathbb{R}^{+N}$ (N -dimensional positive real space). Upon receiving an instance $\vec{X}_t \in [0, 1]^N$, Winnnow makes its prediction $\hat{\ell}_t = +1$ if $W_t = \vec{w}_t \cdot \vec{X}_t \geq \theta$ and -1 otherwise ($\theta > 0$ is a threshold). Given the true label ℓ_t , the weights are updated as follows: $w_{t+1,i} = w_{t,i} \alpha^{X_{t,i}(\ell_t - \hat{\ell}_t)/2}$ for some $\alpha > 1$. If $w_{t+1,i} > w_{t,i}$ we call the update a *promotion* and if $w_{t+1,i} < w_{t,i}$ we call it a *demotion*. Littlestone [21] showed that if each example is labeled by some monotone disjunction of K of its N inputs, then Winnnow will never make more than $O(K \log N)$ prediction mistakes on any sequence of examples. This makes Winnnow a natural tool to apply to learning DNF since by enumerating all 3^n possible terms as inputs to Winnnow, K -term DNF can be learned while making only $O(Kn)$ prediction mistakes. However, the time complexity of running Winnnow this way is exponential in n .

Similar to Winnnow, the Perceptron algorithm maintains a weight vector $\vec{w} \in \mathbb{R}^N$. Upon receiving an instance $\vec{X}_t \in [0, 1]^N$, it makes its prediction $\hat{\ell}_t = +1$ if $W_t = \vec{w}_t \cdot \vec{X}_t \geq \theta$ and -1 otherwise.⁸ Given the true label ℓ_t , the weights are updated as follows: $w_{t+1,i} = w_{t,i} + \alpha X_{t,i}(\ell_t - \hat{\ell}_t)/2$ for some $\alpha > 0$. In contrast to Winnnow, the Perceptron algorithm can be forced to make $\Omega(KN)$ mistakes on monotone K -disjunctions over N inputs [20], making it inappropriate for learning DNF (see also Khardon et al. [18]). However, the additive nature of the weight updates yields much better time complexity bounds for MCMC in contrast to those for multiplicative weight-update schemes (Section 4.1.2).

Inputs to the Weighted Majority algorithm [24] are themselves predictions of “experts” on the current example⁹ \vec{x}_t . Each such expert e_i in the pool has its own weight w_i (initialized to 1), and when a new example \vec{x}_t is given to each expert in the pool, expert e_i sends to WM its prediction $X_{t,i} = e_i(\vec{x}_t) \in \mathbb{R}$, where the sign of $X_{t,i}$ indicates e_i ’s predicted label and $|X_{t,i}|$ can be thought of as e_i ’s confidence (though some experts may restrict themselves to predictions from $\{-1, +1\}$). WM then takes a weighted combination of the predictions and predicts $\hat{\ell}_t = +1$ if $W_t = \vec{w}_t \cdot \vec{X}_t \geq 0$ and -1 otherwise. Upon receiving the correct label ℓ_t , if WM makes a prediction mistake, it reduces the weights of all experts that predicted incorrectly by dividing them by some constant $\alpha > 1$. It has been shown that if the best expert in the pool makes at most v mistakes, then WM has a mistake bound¹⁰ of $O(v + \log N)$. Applying this to predicting nearly as well as the best pruning of an ensemble is straightforward. By placing each possible pruning into the pool, we get a pool size of $N = 2^n$ and thus a mistake bound of $O(v + n)$. However, the time complexity of a straightforward implementation of this algorithm is exponential in n .

2.2. Learning DNF formulas

Learning DNF formulas has been heavily studied, but positive learning-theoretic results exist only in restricted cases, including assuming a uniform distribution over examples (e.g. [5]) or

⁸For additive weight update algorithms like Perceptron, often the threshold is included in the weight vector as $w_{t,0}$, corresponding to an extra attribute $X_{t,0} = 1$. The dot product is then compared to 0 rather than θ .

⁹Throughout this paper, lower case \vec{x} and \vec{y} will represent examples in the original space, while capital \vec{X} and \vec{Y} represent the examples mapped to a new space, which is the input space of Winnnow, Perceptron, and WM.

¹⁰Stronger results on predicting with expert advice were given by Cesa-Bianchi et al. [6] using a more complex algorithm, but these are only better than WM’s by a constant factor. Thus for simplicity, we use WM.

assuming that the number of terms is bounded by $O(\log n)$, where n is the number of variables [4]. In both of these cases, the algorithms require, in addition to labeled examples, *membership queries*, i.e. they need to be able to present arbitrary examples to an oracle and be told their labels.

In contrast, directly applying Winnow to this problem by enumerating all possible terms and learning a monotone disjunction over them does not require any restrictions or the use of membership queries. Since there are only 3^n possible DNF terms over n variables, Winnow's mistake bound on this problem is $O(Kn)$, where K is the number of relevant terms in the target function. However, the time complexity to make a prediction on each example is exponential in this case if a brute-force approach is taken. Indeed, Khardon et al. [18] showed that if $P \neq \#P$, then there is no polynomial time algorithm to exactly simulate Winnow over exponentially many conjunctive features for learning even monotone DNF. Further, while they did provide a kernel allowing them to exactly compute Perceptron's weighted sums when learning DNF, they also gave an exponential lower bound on the number of mistakes that kernel perceptron makes in learning DNF: $2^{\Omega(n)}$.

2.3. Pruning ensembles of classifiers

Prior work in pruning ensembles of classifiers [26,40] (produced by boosting) has been conducted for two reasons. First, the time required to evaluate a complete ensemble is prohibitive in some applications. Second, despite some evidence to the contrary [32,34], boosting can be prone to overfitting. The methods of Margineantu and Dietterich [26] and Tamon and Xiang [40] not only sought subsets of the ensemble with high prediction accuracy, but also with high *diversity*, i.e. hypotheses with high accuracy on different portions of the instance space. The approaches they used included simple ones like early stopping, ones that utilized divergence measures such as Kullback–Leibler divergence or the κ statistic, and methods that used prediction error, sometimes combined with a divergence measure.

To address the concern of overfitting, one can use the WM algorithm, using all possible prunings as “experts” in a pool. Since WM is guaranteed to not perform much worse (in terms of number of on-line prediction mistakes) than the best expert in the pool, we know that a brute-force implementation of this algorithm is guaranteed to not perform much worse than the best pruning. However, a brute-force implementation of WM would take time exponential in the number of hypotheses in the ensemble. So we use an MCMC approach to approximate the weighted sum of the experts' predictions.

2.4. Markov chain Monte Carlo methods

MCMC methods [16] have been applied to problems in combinatorial optimization and approximate summation, where the goal is to approximate weighted sum $W = \sum_{x \in \Omega} s(x)$, where s is a positive function defined on Ω and Ω is a very large, finite set of combinatorial structures. The process involves defining an ergodic Markov chain \mathcal{M} with state space Ω and stationary distribution π . Then one repeatedly simulates \mathcal{M} some number of steps to draw several samples almost according to π . Under appropriate conditions, this technique yields accuracy guarantees. E.g. in approximate summation, sometimes one can guarantee that the estimate of the sum is within a factor ε of the true value with high probability. When this is true and the estimation

algorithm requires only polynomial time, the algorithm is called a *fully polynomial randomized approximation scheme* (FPRAS). In certain cases a similar argument can be made about combinatorial optimization problems, i.e. that the algorithm's solution is within a factor of ε of the true maximum or minimum.

A well-studied problem with an MCMC solution is the *approximate knapsack problem*, where one is given a positive real vector \vec{w} and real number b . The goal is to estimate $|\Omega|$ within a multiplicative factor of ε , where $\Omega = \{\vec{p} \in \{0, 1\}^n : \vec{w} \cdot \vec{p} \leq b\}$ (i.e. as an approximate summation problem, $s(\vec{p}) = 1$ for all $\vec{p} \in \Omega$). Dyer et al. [10] gave a Markov chain for this problem and argued that a polynomial (in n and $1/\varepsilon$) number of samples from it were sufficient to accurately estimate $|\Omega|$. Later, Morris and Sinclair [28] showed that it is sufficient to simulate the chain for a polynomial number of steps to obtain each sample (i.e. that the chain is *rapidly mixing*), thus giving a FPRAS for the knapsack problem.

Another problem with a FPRAS [16] is computing the sum of the weights of a *weighted matching* with parameter λ . For a graph G and $\lambda \geq 0$, approximate $Z_G(\lambda) = \sum_{k=0}^n m_k \lambda^k$, where m_k is the number of matchings in G of size k and n is the number of nodes. This problem has applications to the monomer-dimer problem of statistical physics. In the next section, we combine the knapsack solution with the matching solution to approximate the weighted sums of inputs of linear learning algorithms.

3. Our general algorithm and Markov chain

In general, our state space Ω will consist of the set of inputs to the learning algorithm under consideration (Perceptron, Winnow, or WM). As such, we can think of the states of Ω as functions that map from examples in the original space (the \vec{x} variables) to the input space of the learning algorithm (the \vec{X} variables). So for a state $\vec{p} \in \Omega$ and an input example \vec{x} , we let $p(\vec{x})$ denote \vec{p} evaluated at \vec{x} . E.g. when learning DNF with binary $p(\vec{x})$ (Section 4.1.1), \vec{p} is a term, \vec{x} is an assignment to the variables, and $p(\vec{x}) = 1$ if \vec{x} satisfies \vec{p} and 0 otherwise.

Depending on the application and choice of $p(\vec{x})$ in Section 4, Ω will take on different forms. When learning DNF with Perceptron or Winnow and binary $p(\vec{x})$, we use $\Omega = \{0, 1\}^n$. When learning DNF with Perceptron or Winnow and $p(\vec{x})$ a linear or logistic function (see Section 4.1.1), we use $\Omega = \prod_{i=1}^n \{0, \dots, k_i\}$ for some integers $k_1, \dots, k_n > 0$ (as described in Section 4.1, k_i is the number of values for feature i in a general DNF representation). When we use WM to prune ensembles, we define Ω^+ (similarly, Ω^-) as the set of prunings that predict +1 (similarly, -1) on the current example. In Section 4.2 we show that Ω^+ and Ω^- are each simply $\{0, 1\}^n$ truncated by a single hyperplane. Since in this case the state space of our Markov chain is truncated, we must take care to not exit the state space during a transition. Hence the need for Step 3 in our algorithm below.

Consider a vector $\vec{p} = (p_1, \dots, p_i, \dots, p_n)$. We say that vector \vec{p}' is a *neighbor* of \vec{p} if and only if \vec{p} and \vec{p}' differ in at most one position. I.e. if and only if $\vec{p}' = (p_1, \dots, p'_i, \dots, p_n)$, where p'_i may or may not equal p_i (if $p'_i = p_i$ then the edge from \vec{p} to \vec{p}' is a self-loop). (Note that if Ω is a truncated hypercube, then \vec{p}' might not be in Ω , even if $\vec{p} \in \Omega$. This is why we test for membership in Step 3

below.) We now define \mathcal{M} as a Markov chain with state space Ω that makes transitions from state $\vec{p} \in \Omega$ to state $\vec{q} \in \Omega$ by the following rules.

- (1) With probability $1/2$ let $\vec{q} = \vec{p}$. Otherwise:
- (2) Let \vec{p}' be a neighbor of \vec{p} selected uniformly at random
- (3) If $\vec{p}' \in \Omega$, then let $\vec{p}'' = \vec{p}'$, else let $\vec{p}'' = \vec{p}$.
- (4) With probability $\min\{1, p''(\vec{x}) w_{\vec{p}''} / (p(\vec{x}) w_{\vec{p}})\}$, let $\vec{q} = \vec{p}''$, else let $\vec{q} = \vec{p}$. Here $w_{\vec{p}}$ is the weight of node \vec{p} in the learning algorithm.

Thus \mathcal{M} is a random walk where the transition probabilities favor nodes with higher weights.

Lemma 1. *If every state in Ω can be reached from every other state, then \mathcal{M} is ergodic with stationary distribution*

$$\pi_t(\vec{p}) = \frac{p(\vec{x}_t) w_{\vec{p}}}{W_t},$$

where $W_t = \sum_{\vec{p} \in \Omega} p(\vec{x}_t) w_{\vec{p}}$, i.e. the weighted sum of inputs over all states (inputs) in Ω when example \vec{x}_t is the current example.

Proof. Since all states in Ω can communicate, \mathcal{M} is irreducible. Also, the self-loop of step 1 ensures aperiodicity. Finally, \mathcal{M} is reversible since the transition probabilities

$$P(\vec{p}, \vec{q}) = \frac{\min\{1, q(\vec{x}_t) w_{\vec{q}} / (p(\vec{x}_t) w_{\vec{p}})\}}{2n} = \frac{\min\{1, \pi_t(\vec{q}) / \pi_t(\vec{p})\}}{2n}$$

(here n is the number of neighbors) satisfy the detailed balance condition $\pi_t(\vec{p})P(\vec{p}, \vec{q}) = \pi_t(\vec{q})P(\vec{q}, \vec{p})$. So \mathcal{M} is ergodic with the stated stationary distribution. \square

For each new trial in an on-line algorithm, the weighted sums we estimate are potentially different. Thus we must conduct a new estimation procedure (with a new Markov chain) for each trial. To simplify notation, for the rest of this paper we will let the index t of each trial be implicit, omitting any subscripts unless necessary. Further, in each trial our algorithm defines multiple Markov chains, each assuming that the weight updates of previous trials were made using a different¹¹ learning rate α_i . Hence the weight $w_{\vec{p}}$ of a node \vec{p} (and hence the stationary distribution π and the sum of weights W) will be functions of both α_i and t , so we use the subscript of α_i to denote these differences, leaving the t implicit.

Recalling the definition of Winnow in Section 2.1, if the initial weight vector is the all 1s vector, the weight of term \vec{p} is $w_{\vec{p}} = \alpha^{\vec{p}}$, where $z_{\vec{p}} = \sum_{\vec{x} \in M} \ell_{\vec{x}} p(\vec{x})$, M is the set of examples for which a

¹¹ Note that, however, the actual sequence of updates made will be the same regardless of α_i . This sequence of updates is determined by running the learning algorithm with the original learning rate α .

prediction mistake is made and $\ell_{\vec{x}} \in \{-1, +1\}$ is example \vec{x} 's label. We will refer to $z_{\vec{p}}$ as node \vec{p} 's *total update*, since from it we can directly compute node \vec{p} 's weight. In fact, if the Perceptron algorithm is started with all weights equal to 0, then node \vec{p} 's weight is $w_{\vec{p}} = \alpha z_{\vec{p}}$. Similarly, starting WM with all weights equal to 1 implies $w_{\vec{p}} = \alpha z_{\vec{p}}$, just like Winnow.

Let B be a bound (over all nodes in Ω) on the magnitudes of the total updates up to the current trial, i.e. $B \geq \max_{\vec{p} \in \Omega} \{|z_{\vec{p}}|\}$. Since this requires taking the maximum over an exponentially large set, we note that it suffices to instead use $B \geq \sum_{\vec{x} \in \mathcal{M}} \max_{\vec{p} \in \Omega} \{|\ell_{\vec{x}} p(\vec{x})|\}$. This quantity is easy to bound so long as bounds on the possible values of $\ell_{\vec{x}}$ and $p(\vec{x})$ are known for each \vec{x} , which is the case for all our algorithms. (E.g. in Winnow and Perceptron, it suffices to set B equal to the sum of all promotions and demotions made on all examples for which a prediction mistake was made up to the current trial.) Now let r be the smallest integer such that $(1 + 1/B)^{r-1} \geq \alpha$ and $r \geq 1 + \log_2 \alpha$ (so $r \leq 2 + B \ln \alpha$). Also, let $\zeta = 1/(\alpha^{1/(r-1)} - 1) \geq B$ and $\alpha_i = (1 + 1/\zeta)^{i-1} = \alpha^{(i-1)/(r-1)}$ for $1 \leq i \leq r$ (so $\alpha_r = \alpha$).

Now define $f_i(\vec{p}) = w_{\alpha_{i-1}\vec{p}}/w_{\alpha_i\vec{p}}$, where \vec{p} is chosen according to π_{α_i} . Then

$$E[f_i] = \sum_{\vec{p} \in \Omega} \left(\frac{w_{\alpha_{i-1}\vec{p}}}{w_{\alpha_i\vec{p}}} \right) \frac{p(\vec{x}) w_{\alpha_i\vec{p}}}{W(\alpha_i)} = \frac{W(\alpha_{i-1})}{W(\alpha_i)}.$$

So we can estimate $W(\alpha_{i-1})/W(\alpha_i)$ by sampling states \vec{p} from \mathcal{M} and computing the sample mean of the $f_i(\vec{p})$. Note that

$$W(\alpha) = \left(\frac{W(\alpha_r)}{W(\alpha_{r-1})} \right) \left(\frac{W(\alpha_{r-1})}{W(\alpha_{r-2})} \right) \cdots \left(\frac{W(\alpha_2)}{W(\alpha_1)} \right) W(\alpha_1).$$

So for each value $\alpha_2, \dots, \alpha_r$, we run S independent simulations of \mathcal{M} , and let X_i be the sample mean of $w_{\alpha_{i-1}}/w_{\alpha_i}$. Then our estimate¹² is

$$\hat{W}(\alpha) = W(\alpha_1) \prod_{i=2}^r 1/X_i. \quad (1)$$

In order to complete our computation of \hat{W} , we must also compute $W(\alpha_1)$. Due to the definition of α_1 , this is straightforward for learning DNF with the various definitions of $p(\vec{x})$ (Sections 4.1.1 and 4.1.2). For the ensemble pruning problem, $W^+(\alpha_1) = |\Omega^+|$, where Ω^+ is the set of prunings that predict +1 on the input \vec{x} . This cannot be efficiently computed exactly, so we must estimate it with the FPRAS of Morris and Sinclair [28] (Section 4.2).

The following theorem bounds the error of our algorithm's estimates of W . The theorem is based on *variation distance*, which is a distance measure between a Markov chain's simulated and stationary distributions, defined as $\max_{U \subseteq \Omega} |P^\tau(\vec{p}, U) - \pi(U)|$, where $P^\tau(\vec{p}, \cdot)$ is the distribution of a chain's state at simulation step τ given that the simulation started in state $\vec{p} \in \Omega$, and π is the chain's stationary distribution.

Theorem 2. Assume $a \leq f_i, \hat{f}_i \leq b$ for all i , where \hat{f}_i is the same as f_i but with samples drawn according to the distribution yielded by simulating \mathcal{M} . Let the sample size $S = \lceil 130rb/(ae^2) \rceil$ and \mathcal{M} be

¹²When we apply our results to the Perceptron algorithm in Section 4.1.2, we will also use $\alpha_0 = 0$ and update the product of ratios accordingly.

simulated long enough for each sample such that the variation distance between the empirical distribution and π_{α_i} is at most $\varepsilon a/(5br)$ for each i . Also, assume that $W(\alpha_1)$ can be computed exactly. Then for any $\delta > 0$, $\hat{W}(\alpha)$ satisfies

$$\Pr[(1 - \varepsilon)W(\alpha) \leq \hat{W}(\alpha) \leq (1 + \varepsilon)W(\alpha)] \geq 1 - \delta.$$

Proof. Let the distribution $\hat{\pi}_{\alpha_i}$ be the one resulting from simulating \mathcal{M} , and assume that the variation distance $\|\hat{\pi}_{\alpha_i} - \pi_{\alpha_i}\| \leq \varepsilon a/(5br)$. Now consider the random variable \hat{f}_i , which is the same as f_i except that the terms are selected according to $\hat{\pi}_{\alpha_i}$. Since $\hat{f}_i \in [a, b]$, $|\mathbb{E}[\hat{f}_i] - \mathbb{E}[f_i]| \leq \varepsilon a/(5r)$, which implies $\mathbb{E}[f_i] - \varepsilon a/(5r) \leq \mathbb{E}[\hat{f}_i] \leq \mathbb{E}[f_i] + \varepsilon a/(5r)$. Factoring out $\mathbb{E}[f_i]$ from both sides and noting that $1/\mathbb{E}[f_i] \leq 1/a$ yields

$$\left(1 - \frac{\varepsilon}{5r}\right)\mathbb{E}[f_i] \leq \mathbb{E}[\hat{f}_i] \leq \left(1 + \frac{\varepsilon}{5r}\right)\mathbb{E}[f_i]. \quad (2)$$

This allows us to conclude that $\mathbb{E}[\hat{f}_i] \geq \mathbb{E}[f_i]/2$. Since $\hat{f}_i \leq b$, we get $\text{Var}[\hat{f}_i] \leq b \mathbb{E}[\hat{f}_i]$, yielding

$$\frac{\text{Var}[\hat{f}_i]}{(\mathbb{E}[\hat{f}_i])^2} \leq \frac{b}{\mathbb{E}[\hat{f}_i]} \leq \frac{2b}{\mathbb{E}[f_i]} \leq 2b/a. \quad (3)$$

Let $X_i^{(1)}, \dots, X_i^{(S)}$ be a sequence of S independent copies of \hat{f}_i , and let $\bar{X}_i = (\sum_{j=1}^S X_i^{(j)})/S$. Then $\mathbb{E}[\bar{X}_i] = \mathbb{E}[\hat{f}_i]$ and $\text{Var}[\bar{X}_i] = \text{Var}[\hat{f}_i]/S$. The estimator of $W(\alpha)$ is $W(\alpha_1)/X = W(\alpha_1)/\prod_{i=2}^r \bar{X}_i$. Since the \bar{X}_i 's are independent, $\mathbb{E}[X] = \prod_{i=2}^r \mathbb{E}[\bar{X}_i] = \prod_{i=2}^r \mathbb{E}[\hat{f}_i]$ and $\mathbb{E}[X^2] = \prod_{i=2}^r \mathbb{E}[\bar{X}_i^2]$. Let $\rho = \prod_{i=2}^r W(\alpha_{i-1})/W(\alpha_i)$, (i.e. what we are estimating with X) and $\hat{\rho} = \mathbb{E}[X]$. Then applying Eq. (2) gives

$$\left(1 - \frac{\varepsilon}{5r}\right)^r \rho \leq \hat{\rho} \leq \left(1 + \frac{\varepsilon}{5r}\right)^r.$$

Since $\lim_{r \rightarrow \infty} (1 + \varepsilon/(5r))^r = e^{\varepsilon/5} \leq 1 + \varepsilon/4$ and $(1 - \varepsilon/(5r))^r$ is minimized at $r = 1$, we get

$$\left(1 - \frac{\varepsilon}{4}\right)\rho \leq \hat{\rho} \leq \left(1 + \frac{\varepsilon}{4}\right)\rho.$$

Since $\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$, we have

$$\begin{aligned} \frac{\text{Var}[X]}{(\mathbb{E}[X])^2} &= \prod_{i=1}^{r-1} \left(1 + \frac{\text{Var}[\bar{X}_i]}{(\mathbb{E}[\bar{X}_i])^2}\right) - 1 \\ &\leq \left(1 + \frac{2b}{aS}\right)^{r-1} - 1 \quad (\text{by Eq.(3)}) \\ &\leq \left(1 + \frac{\varepsilon^2}{65r}\right)^r - 1 \leq \exp(\varepsilon^2/65) - 1 \leq \varepsilon^2/64. \end{aligned}$$

The last inequality holds since $\exp(x/65) \leq 1 + 1/64$ for $x \in [0, 1]$. We now apply Chebyshev's inequality to X with standard deviation $\varepsilon\hat{\rho}/8$:

$$\Pr[|X - \hat{\rho}| > \varepsilon\hat{\rho}/4] \leq 1/4.$$

So with probability at least $3/4$ we get

$$\left(1 - \frac{\varepsilon}{4}\right)\hat{\rho} \leq X \leq \left(1 + \frac{\varepsilon}{4}\right)\hat{\rho},$$

which implies that with probability at least $3/4$

$$\frac{(1 + \varepsilon)}{\rho} \geq \frac{1}{(1 - \varepsilon/4)^2 \rho} \geq 1/X \geq \frac{1}{(1 + \varepsilon/4)^2 \rho} \geq \frac{(1 - \varepsilon)}{\rho}. \quad (4)$$

Making the approximation with probability at least $1 - \delta$ for any $\delta > 0$ is done by rerunning $O(\ln 1/\delta)$ times the procedure for estimating X and taking the median of the results [17]. \square

It is also possible to extend Theorem 2 to the case where $W(\alpha_1)$ cannot be exactly computed, but can be accurately estimated.

Corollary 3. Assume $a \leq f_i, \hat{f}_i \leq b$ for all i . Let the sample size $S = \lceil 30 r b / (a \varepsilon^2) \rceil$, $W(\alpha_1)$'s estimate be within $\varepsilon/2$ of its true value with probability $\geq 3/4$, and \mathcal{M} be simulated long enough for each sample such that the variation distance between the empirical distribution and π_{α_i} is at most $\varepsilon a / (10br)$ for all i . Then for any $\delta > 0$, $\hat{W}(\alpha)$ satisfies

$$\Pr[(1 - \varepsilon)W(\alpha) \leq \hat{W}(\alpha) \leq (1 + \varepsilon)W(\alpha)] \geq 1 - \delta.$$

Proof. The analysis is the same as in the proof of Theorem 2, except we now must accommodate another source of error. First, substitute $\varepsilon/2$ for ε in Eq. (4). Given the accuracy of $\hat{W}(\alpha_1)$ with probability at least $3/4$, we get

$$\frac{W(\alpha_1)(1 - \varepsilon/2)^2}{\rho} \leq \frac{\hat{W}(\alpha_1)}{X} \leq \frac{W(\alpha_1)(1 + \varepsilon/2)^2}{\rho},$$

with probability at least $1/2$. This completes the proof (the constants in S remain unchanged). Similar to Theorem 2, both estimates can be run multiple times and the median taken in order to reduce the probability of failure. \square

We now bound the mixing time of \mathcal{M} by using the *canonical paths* method [38]. In this method, we treat \mathcal{M} as a directed graph with vertices Ω and edges $E = \{(\vec{p}, \vec{q}) \in \Omega \times \Omega: Q(\vec{p}, \vec{q}) > 0\}$, where $Q(\vec{p}, \vec{q}) = \pi_{\alpha}(\vec{p}) P(\vec{p}, \vec{q})$. For each ordered pair $(\vec{p}, \vec{q}) \in \Omega \times \Omega$, we specify a canonical path $\gamma_{\vec{p}, \vec{q}} \in \Gamma$ from \vec{p} to \vec{q} in the graph (Ω, E) that corresponds to a set of legal transitions in \mathcal{M} from \vec{p} to \vec{q} . We measure how heavily any one edge in E is loaded with canonical paths by

$$\bar{\rho} = \bar{\rho}(\Gamma) = \max_{e \in E} \left\{ \frac{1}{Q(e)} \sum_{\gamma_{\vec{p}, \vec{q}} \ni e} \pi_{\alpha}(\vec{p}) \pi_{\alpha}(\vec{q}) |\gamma_{\vec{p}, \vec{q}}| \right\}. \quad (5)$$

We start with a result from Sinclair [38], restated by Jerrum and Sinclair [16].

Theorem 4 (Jerrum and Sinclair [16], Sinclair [38]). *Let \mathcal{M} be a finite, reversible, ergodic Markov chain with loop probabilities $P(\vec{p}, \vec{p}) \geq 1/2$ for all \vec{p} . Let Γ be a set of canonical paths with maximum edge loading $\bar{\rho} = \bar{\rho}(\Gamma)$. Then the mixing time of \mathcal{M} satisfies $\tau_{\vec{p}}(\varepsilon) \leq \bar{\rho}(\ln 1/\pi(\vec{p}) + \ln 1/\varepsilon)$ for any choice of initial state \vec{p} , i.e. after simulating \mathcal{M} for $\bar{\rho}(\ln 1/\pi(\vec{p}) + \ln 1/\varepsilon)$ steps starting in \vec{p} , the variation distance between $\hat{\pi}_{\alpha_i}$ and π_{α_i} is at most ε .*

In general, $\Omega = \prod_{i=1}^n \{0, \dots, k_i\}$ for some integers k_1, \dots, k_n . Without loss of generality we let $\Omega = \{0, \dots, k\}^n$ for some positive integer k . Then there is an edge from node $\vec{p} = (p_1, \dots, p_i, \dots, p_n)$ to $\vec{p}' = (p_1, \dots, p'_i, \dots, p_n)$, i.e. an edge exists between each pair of nodes that differ in at most one position (self-loops also exist). For our proof, we assume that the hypercube is untruncated, which is necessary to ensure that no canonical paths leave the chain. However, it is likely that mixing time bounds also exist for truncated hypercubes. Such a bound could probably be derived by the recent work of Morris and Sinclair [28] who give an FPRAS for a truncated Boolean hypercube that has a uniform distribution.

Let $\vec{p} = (p_1, \dots, p_n)$ and $\vec{q} = (q_1, \dots, q_n)$ be arbitrary states of Ω . The canonical path $\gamma_{\vec{p}, \vec{q}}$ consists of n edges, where edge i is

$$((q_1, \dots, q_{i-1}, p_i, p_{i+1}, \dots, p_n), (q_1, \dots, q_{i-1}, q_i, p_{i+1}, \dots, p_n)),$$

i.e. position i is changed from \vec{p}_i to \vec{q}_i . So some edges of $\gamma_{\vec{p}, \vec{q}}$ might be loops. Now focus on a particular oriented edge

$$e = (\vec{a}, \vec{a}') = ((a_1, \dots, a_i, \dots, a_n), (a_1, \dots, a'_i, \dots, a_n)).$$

We will now bound Eq. (5) for e , which yields a bound on $\bar{\rho}$ and allows us to apply Theorem 4. Let $\text{cp}(e) = \{(\vec{p}, \vec{q}) : \gamma_{\vec{p}, \vec{q}} \ni e\}$ be the set of endpoints of canonical paths that use edge e . We use Jerrum and Sinclair's [16] mapping $\vec{\eta}_e : \text{cp}(e) \rightarrow \Omega$, defined¹³ as follows: if $(\vec{p}, \vec{q}) = ((p_1, \dots, p_n), (q_1, \dots, q_n)) \in \text{cp}(e)$, then

$$\vec{\eta}_e(\vec{p}, \vec{q}) = (b_1, \dots, b_n) = (p_1, \dots, p_{i-1}, a_i, q_{i+1}, \dots, q_n).$$

Note that $\vec{p} = (b_1, \dots, b_{i-1}, a_i, a_{i+1}, \dots, a_n)$ and $\vec{q} = (a_1, \dots, a_{i-1}, a'_i, b_{i+1}, \dots, b_n)$. Since \vec{p} and \vec{q} can be unambiguously recovered from $\vec{\eta}_e(\vec{p}, \vec{q})$, the mapping $\vec{\eta}_e$ is injective.

We are now ready to state the mixing time bound.

Theorem 5. *For all $\vec{p}, \vec{q} \in \Omega$ and for all $e \in \Omega \times \Omega$ such that $(\vec{p}, \vec{q}) \in \text{cp}(e)$, assume*

$$\pi_{\alpha}(\vec{p}) \pi_{\alpha}(\vec{q}) \leq g \pi_{\alpha}(\vec{a}') \pi_{\alpha}(\vec{\eta}_e(\vec{p}, \vec{q}))$$

for some function $g = g(n, K, k, \alpha)$. Also assume that for all neighbors \vec{a} and \vec{a}' in Ω ,

$$\max\{\pi_{\alpha}(\vec{a})/\pi_{\alpha}(\vec{a}'), \pi_{\alpha}(\vec{a}')/\pi_{\alpha}(\vec{a})\} \leq h = h(n, K, k, \alpha).$$

Then a simulation of \mathcal{M} that starts at node \vec{p} and is of length

$$T = 2kn^2 g h \left(\ln \left(\frac{W(\alpha)}{w_{\alpha, \vec{p}}} \right) + \ln(1/\varepsilon') \right)$$

will draw samples from $\hat{\pi}_{\alpha}$ such that $\|\hat{\pi}_{\alpha} - \pi_{\alpha}\| \leq \varepsilon'$.

¹³ Vector notation is used when denoting $\vec{\eta}_e$ since $\vec{\eta}_e(\vec{p}, \vec{q}) \in \Omega$ for all $\vec{p}, \vec{q} \in \Omega$.

Proof. Since $Q(e) = \min\{\pi_\alpha(\vec{a}), \pi_\alpha(\vec{a}')\}/(2kn)$, we get

$$\begin{aligned}\pi_\alpha(\vec{p}) \pi_\alpha(\vec{q}) &\leq \frac{2kn g Q(e)}{\min\{\pi_\alpha(\vec{a}), \pi_\alpha(\vec{a}')\}} \pi_\alpha(\vec{a}') \pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q})) \\ &= 2kn g Q(e) \max\{1, \pi_\alpha(\vec{a}')/\pi_\alpha(\vec{a})\} \pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q})) \\ &\leq 2kn g h Q(e) \pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q})).\end{aligned}$$

Given the above inequality, we can now bound $\bar{\rho}$. Since $|\gamma_{\vec{p}, \vec{q}}| = n$, we get

$$\frac{1}{Q(e)} \sum_{\gamma_{\vec{p}, \vec{q}} \ni e} \pi_\alpha(\vec{p}) \pi_\alpha(\vec{q}) |\gamma_{\vec{p}, \vec{q}}| \leq 2kn^2 gh \sum_{\gamma_{\vec{p}, \vec{q}} \ni e} \pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q})) \leq 2kn^2 gh.$$

The last inequality holds because $\vec{\eta}_e$ is injective and π_α is a probability distribution. Applying Theorem 4 completes the proof. \square

Corollary 6. *For Markov chains for which g and h are polynomial in n and K (we assume k and α are constants) and for approximation schemes for which b and $1/a$ are polynomial in n and K , our algorithm is an FPRAS.*

4. Example applications

4.1. Learning DNF formulas

We consider generalized DNF representations, where the instance space is $\prod_{i=0}^{n-1} \{1, \dots, k_i\}$ and the set of terms is $\prod_{i=0}^{n-1} \{0, \dots, k_i\}$, where k_i is the number of values for feature i . A term $\vec{p} = (p_0, \dots, p_{n-1})$ is satisfied by example $\vec{x} = (x_0, \dots, x_{n-1})$ if and only if $\forall p_i > 0, p_i = x_i$. So $p_i = 0$ implies that x_i is irrelevant for term \vec{p} , and $p_i > 0$ implies that x_i must equal p_i for \vec{p} to be satisfied.

We present algorithms to learn this concept class that are based on Littlestone's Winnow [21] and Rosenblatt's Perceptron [33] algorithms. The inputs to the linear threshold units learned by these algorithms consist of the entire set of DNF terms over the original set of n inputs. For reasons that will become clear, we look at different versions of the function $p(\vec{x})$, which measures the degree to which \vec{x} satisfies \vec{p} . These versions include $p(\vec{x})$ being a threshold function, a logistic function, and a linear function.

None of our approaches in this section give complete, efficient solutions to the problem of learning DNF in the on-line model, but they do give new mechanisms that could be refined for potential application to restricted subclasses of learning DNF, e.g. restricted classes of DNF or specific distributions over the examples.

4.1.1. Winnow

Recalling the definition of Winnow in Section 2.1, if the initial weight vector is the all 1s vector, the weight of term \vec{p} is $w_{\vec{p}} = \alpha^{z_{\vec{p}}}$, where $z_{\vec{p}} = \sum_{\vec{x} \in M} \ell_{\vec{x}} p(\vec{x})$, M is the set of examples for which a

prediction mistake is made and $\ell_{\vec{x}} \in \{-1, +1\}$ is example \vec{x} 's label. We now state one of Littlestone's results for Winnow [23], which we will use to bound the number of prediction mistakes it makes for the variations of our algorithm.

Theorem 7 (Littlestone [23]). *Let $(\vec{Y}_j, \ell_j) \in [0, 1]^N \times \{0, 1\}$ for $j = 1, \dots, t$ (t is the index of the current trial). Suppose that there exist $\vec{\mu} \geq 0$ and $0 < \rho < 1$ such that whenever $\ell_j = 1$ we have $\vec{\mu} \cdot \vec{Y}_j \geq 1$ and whenever $\ell_j = 0$ we have $\vec{\mu} \cdot \vec{Y}_j \leq 1 - \rho$. Now suppose Winnow sees as inputs $\mathcal{X} = (\vec{X}_j, \ell_j)$ where each $X_{ij} \in [0, 1]$, and define $\vec{E}_j = (|X_{j1} - Y_{j1}|, \dots, |X_{jN} - Y_{jN}|)$. Then the number of mistakes made by Winnow on \mathcal{X} with $\alpha = 1 + \rho/2$ and $\theta = N$ is at most*

$$8/\rho^2 + \max\left(0, \frac{14}{\rho^2} \sum_{i=1}^N \mu_i \ln(\mu_i \theta)\right) + \frac{4}{\rho} \sum_{j=1}^t \vec{\mu} \cdot \vec{E}_j.$$

We will examine three versions of our algorithm, differing in the values that Winnow receives as its inputs. In the following, we say that a variable x_i in example \vec{x} *matches* its corresponding variable p_i in term \vec{p} if $p_i = 0$ or $p_i = x_i$. We let $m_{\vec{x}, \vec{p}} \in \{0, \dots, n\}$ denote the number of variables in \vec{x} that match their corresponding variables in \vec{p} (we drop the subscript \vec{x} when it is clear from context).

- (1) *Binary* $p(\vec{x})$ means that Winnow input $X_{\vec{p}} = 1$ if $m_{\vec{p}} = n$ and 0 otherwise.
- (2) *Logistic* $p(\vec{x})$ means that Winnow input $X_{\vec{p}}$ is

$$p(\vec{x}) = \frac{2}{1 + e^{-\sigma(m_{\vec{p}} - n)}}, \quad (6)$$

where $\sigma > 0$ is a parameter. Thus $p(\vec{x}) \in (0, 1]$ and grows as \vec{p} becomes more satisfied by \vec{x} (it equals 1 if and only if it is completely satisfied).

- (3) *Linear* $p(\vec{x})$ means that Winnow input $X_{\vec{p}}$ is

$$p(\vec{x}) = \frac{1 + m_{\vec{p}}}{n + 1}. \quad (7)$$

Thus $p(\vec{x}) \in (0, 1]$ and grows as \vec{p} becomes more satisfied by \vec{x} (it equals 1 if and only if it is completely satisfied).

We defined $p(\vec{x}) > 0$ for all \vec{x} in order to ensure that for a given \vec{x} , every term in Ω contributes something to the weighted sum, and the hypercube is untruncated, which allows us to apply Theorem 5. Using binary $p(\vec{x})$ also yields an untruncated hypercube, as explained below.

Before considering these three cases individually, we state some common results for them all. First note that for logistic and linear $p(\vec{x})$, Ω consists of the entire set of possible terms, since each term gives to Winnow a value $p(\vec{x}) > 0$. Thus in the chain defined in Section 3, every state can be reached from every other state. Further, for binary $p(\vec{x})$, we note that there are exactly 2^n terms that are satisfied by \vec{x} , i.e. \vec{p} is satisfied by \vec{x} if and only if $p_i = 0$ or $p_i = x_i$ for all $i \in \{1, \dots, n\}$. Thus in this case we can construct the state space to be $\Omega = \{0, 1\}^n$, which is completely connected and untruncated. Therefore it is obvious that Lemma 1 applies to all our Markov chains. We now discuss the application of Theorem 2. All that is required to apply this result is to bound the range

of f and \hat{f} . Since f and \hat{f} are independent of $p(\vec{x})$, the same result applies to all three versions of our algorithm.

Lemma 8. *When applying binary, logistic, or linear Winnow to learn DNF, for all i , $1/e \leq f_i, \hat{f}_i \leq e$.*

Proof. First note that the only difference between f_i and \hat{f}_i is the probability distribution that generates the terms that define them, i.e. their ranges are the same. Thus we focus on bounding f_i only. Let $z_{\vec{p}}$ be node \vec{p} 's total update as defined in Section 3. Then

$$f_i(\vec{p}) = \frac{w_{\alpha_{i-1}, \vec{p}}}{w_{\alpha_i, \vec{p}}} = \left(\frac{\alpha_{i-1}}{\alpha_i} \right)^{z_{\vec{p}}} = \left(\frac{(1 + 1/\zeta)^{i-2}}{(1 + 1/\zeta)^{i-1}} \right)^{z_{\vec{p}}} = (1 + 1/\zeta)^{-z_{\vec{p}}}.$$

Recall that from its definition, $\zeta \geq B \geq |z_{\vec{p}}|$ for all \vec{p} (to avoid division by zero, we can also assume that $\zeta > 0$). If $z_{\vec{p}} < 0$, then $1 \leq (1 + 1/\zeta)^{-z_{\vec{p}}} \leq e$. If $z_{\vec{p}} \geq 1$, then $1/e \leq (1 + 1/\zeta)^{-z_{\vec{p}}} \leq 1$. \square

Note that $W(\alpha_1) = W(1)$ is simply $\sum_{\vec{p} \in \Omega} p(\vec{x})$. For the binary case, this is simply the number of terms satisfied by \vec{x} , which equals 2^n . For the linear and logistic cases, it can be efficiently computed exactly if we assume that $k_i = k$ for all i . Under this assumption, the number of terms that match exactly $i \in \{0, \dots, n\}$ variables in the example \vec{x} is $2^i \binom{n}{i} (k-1)^{n-i}$, since there are $\binom{n}{i}$ positions to place the matched variables, each matched position p_j can equal 0 or x_j , and each unmatched position $p_{j'}$ can take on any value from $\{1, \dots, k\} \setminus \{x_{j'}\}$. Thus for the logistic case, we get

$$W(1) = \sum_{i=0}^n 2^i \binom{n}{i} \frac{2(k-1)^{n-i}}{1 + e^{-\sigma(i-n)}}, \quad (8)$$

and for the linear case, we get

$$W(1) = \sum_{i=0}^n 2^i \binom{n}{i} \frac{(i+1)(k-1)^{n-i}}{n+1} = \frac{(k+1)^n + 2n(k+1)^{n-1}}{n+1}. \quad (9)$$

Thus all three can efficiently be computed exactly. By applying this and substituting Lemma 8's bounds into Theorem 2, we get the following.

Corollary 9. *When applying Winnow to learn generalized DNF (with $k_i = k$ for all i for the logistic and linear cases), let the sample size $S = \lceil 130 r e^2 / \epsilon^2 \rceil$ and \mathcal{M} be simulated long enough for each sample such that the variation distance between the empirical distribution and $\pi_{\alpha_{i,t}}$ is at most $\epsilon / (5e^2 r)$. Then for any $\delta > 0$, $\hat{W}(\alpha)$ satisfies*

$$\Pr[(1 - \epsilon)W(\alpha) \leq \hat{W}(\alpha) \leq (1 + \epsilon)W(\alpha)] \geq 1 - \delta.$$

As stated in the following corollary, our algorithms' behaviors are the same as Winnow's for a straightforward brute-force implementation if the weighted sums are not too close to θ for any input. This is true with probability at least $1 - \delta'_t$ for each trial t , so setting $\delta'_t = \delta / 2^t$ yields a total

probability of failure of at most¹⁴ $\sum_{t=1}^{\infty} \delta/2^t = \delta$. Finally, note that it is easy to extend the corollary to tolerate a bounded number of trials with weighted sums that are near θ by thinking of such potential mispredictions as noise and applying Theorem 7.

Corollary 10. *Using the assumptions of Corollary 9, if $W_t(\alpha) \notin [\theta/(1+\varepsilon), \theta/(1-\varepsilon)]$ for all trials t , then with probability at least $1-\delta$, the number of mistakes made by Winnow on any sequence of examples is as bounded by Theorem 7 (see Sections 4.1.1.1–4.1.1.3 and Lemma 11).*

A hurdle that must be overcome to get an efficient algorithm is S 's polynomial dependence on $1/\varepsilon$ in Corollary 9, even though Winnow might at times have W/θ exponentially close¹⁵ to 1, requiring exponentially small ε . It is open whether this can be addressed in an average-case analysis of Winnow when learning restricted concept classes under specific distributions.

We now explore bounding the mixing times of the Markov chains. Note that the bounds are based on worst-case analyses and assume that the maximum number of weight updates (as bounded by Theorem 7's mistake bound) have been made. Prior to making that number of updates (e.g. near the start of training), the mixing time bounds will be lower since the stationary distribution π of \mathcal{M} will be closer to uniform (in fact, before the first update, π is uniform). We can get mixing time bounds for these earlier cases by substituting the number of prediction mistakes made so far for the mistake bounds.

4.1.1.1. Binary $p(\vec{x})$. It is straightforward to apply Theorem 7 to the binary case. Let the vector $\vec{\mu}$ be 0 for each irrelevant term and 1 for each relevant term. Then when $\ell = 1$, at least one relevant term must be satisfied, so $\vec{\mu} \cdot \vec{Y} \geq 1$. Further, if $\ell = 0$, then no relevant terms are satisfied and $\vec{\mu} \cdot \vec{Y} = 0 \leq 1 - \rho$ for $\rho = 1$. Assuming all examples are noise-free, applying Theorem 7 yields a mistake bound of $|M| \leq 8 + 14K \ln N$. So if $k \geq k_i$ for all i , then using the at most $(k+1)^n$ possible terms as Winnow's inputs, it can learn K -term generalized DNF with at most $8 + 14Kn \ln(k+1)$ prediction mistakes.

Unfortunately, with the binary case it is very difficult to find non-trivial bounds on g and h from Theorem 5, due to the discontinuity of $p(\vec{x})$. Bounding both g and h requires bounding the ratios of the weights of nodes in Ω . For binary $p(\vec{x})$, these weights directly depend on how often the nodes predicted 1 when a prediction mistake was made, but it is difficult to relate how often this occurs for a node \vec{p} to how often this occurs for another node \vec{q} , even if \vec{p} and \vec{q} are neighbors. On the other hand, when we consider logistic and linear $p(\vec{x})$, we can relate the nodes' weights and get non-trivial bounds on g and h .

4.1.1.2. Logistic $p(\vec{x})$. The mistake bound of this application of Winnow is similar to that of the straightforward version with binary inputs.

¹⁴ Recall from the proof of Theorem 2 that only $O(\log 1/\delta')$ runs of the estimation procedure are needed to reduce the probability of failure to δ' .

¹⁵ The potential problem of $W_i/\theta = 1$ can be avoided by using a threshold of $\theta + \alpha^{-(|M|+1)}$, where $|M|$ is the mistake bound from applying Theorem 7. Obviously W can never equal this new threshold.

Lemma 11. When using Equation 6 with $\sigma = \ln(60Kn \ln k)$ to specify the inputs, where $k = \max_i \{k_i\}$, the number of prediction mistakes made by Winnow when learning DNF is at most

$$8.88 + 15.54Kn \ln k.$$

Proof. We start by finding $\vec{\mu}$ and ρ that satisfy the conditions of Theorem 7. For each of the K relevant terms (Winnow inputs), set the corresponding value in $\vec{\mu}$ to a constant μ , which we will define later. Set all other values in $\vec{\mu}$ to 0. In the worst case, when an example \vec{x} is positive, it satisfies exactly one relevant term \vec{p} and does not at all satisfy any of the other relevant terms. Then $p(\vec{x}) = 1$ and $q(\vec{x}) = 2/(1 + e^{\sigma n})$ for all other relevant terms \vec{q} . Thus it suffices to set μ such that

$$\mu + \frac{2(K-1)\mu}{1 + e^{\sigma n}} \geq 1.$$

After some algebra we see that it suffices to set $\mu = (1 + e^{\sigma n})/(2K + e^{\sigma n} - 1)$.

Now we find ρ . In the worst case, for a negative example \vec{x} , we will have each relevant term \vec{p} almost fully satisfied, i.e. $m_{\vec{p}} = n - 1$. Hence $p(\vec{x}) = 2/(1 + e^{\sigma})$. So we need ρ such that $1 - \rho \geq 2K\mu/(1 + e^{\sigma})$. Substituting μ yields

$$\rho \leq 1 - \frac{2K(1 + e^{\sigma n})}{2K(1 + e^{\sigma}) + e^{\sigma n}(1 + e^{\sigma}) - e^{\sigma} - 1}.$$

This expression decreases with increasing n , so to find an appropriate ρ , it suffices to take its limit as $n \rightarrow \infty$. Applying l'Hôpital's rule shows that $\rho = 1 - 2K/(1 + e^{\sigma})$ is sufficient, which is positive so long as $\sigma > \ln(2K - 1)$. We assume that all examples are noise-free, so $\vec{E}_j = \vec{0}$ for all j . Now applying Theorem 7 yields a mistake bound of

$$\begin{aligned} & \left(\frac{(1 + e^{\sigma})^2}{(1 + e^{\sigma})^2 - 4K(e^{\sigma} + 1 - K)} \right) \\ & \cdot \left(8 + 14K \left(\frac{1 + e^{\sigma n}}{2K + e^{\sigma n} - 1} \right) \left(\ln \left(\frac{1 + e^{\sigma n}}{2K + e^{\sigma n} - 1} \right) + \ln N \right) \right) \\ & \leq \left(\frac{(1 + e^{\sigma})^2}{(1 + e^{\sigma})^2 - 4K(1 + e^{\sigma})} \right) (8 + 14K \ln N) \\ & = \left(\frac{1 + e^{\sigma}}{1 + e^{\sigma} - 4K} \right) (8 + 14K \ln N) \leq 1.11(8 + 14K \ln N), \end{aligned}$$

since $K, n \geq 1$ and $k \geq 2$. Noting that $N \leq k^n$ completes the proof. \square

We now work towards a mixing time bound for the chain.

Lemma 12. Let $\Omega = \{0, \dots, k\}^n$. Then for all $\vec{p}, \vec{q} \in \Omega$,

$$\pi_{\alpha}(\vec{p})\pi_{\alpha}(\vec{q}) \leq 4\alpha^{|M|}\pi_{\alpha}(\vec{a}')\pi_{\alpha}(\vec{\eta}_e(\vec{p}, \vec{q})).$$

Proof. Let $\vec{q}_{1\dots i}$ denote (q_1, \dots, q_i) , and similarly for $\vec{p}_{1\dots i}$. Then $m_{\vec{p}} = m_{\vec{p}_{1\dots i}} + m_{\vec{p}_{i+1\dots n}}$, $m_{\vec{q}} = m_{\vec{q}_{1\dots i}} + m_{\vec{q}_{i+1\dots n}}$, $m_{\vec{\eta}_e(\vec{p}, \vec{q})} = m_{\vec{p}_{1\dots i}} + m_{\vec{q}_{i+1\dots n}}$, and $m_{\vec{a}'} = m_{\vec{q}_{1\dots i}} + m_{\vec{p}_{i+1\dots n}}$. Further, all four of these values are in $\{0, \dots, n\}$. This yields

$$\begin{aligned} \frac{\pi_\alpha(\vec{p})\pi_\alpha(\vec{q})}{\pi_\alpha(\vec{a}')\pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q}))} &= \frac{p(\vec{x})q(\vec{x})\alpha^{z_{\vec{p}}+z_{\vec{q}}}}{\vec{\eta}_e(\vec{p}, \vec{q})(\vec{x})\vec{a}'(\vec{x})\alpha^{z_{\vec{\eta}_e(\vec{p}, \vec{q})}+z_{\vec{a}'}}} \\ &= \left(\frac{\alpha^{z_{\vec{p}}+z_{\vec{q}}}}{\alpha^{z_{\vec{\eta}_e(\vec{p}, \vec{q})}+z_{\vec{a}'}}} \right) \left(\frac{1 + U(\vec{p}_{1\dots i}, \vec{q}_{i+1\dots n}, n) + U(\vec{q}_{1\dots i}, \vec{p}_{i+1\dots n}, n) + U(\vec{p}_{1\dots n}, \vec{q}_{1\dots n}, 2n)}{1 + U(\vec{p}_{1\dots i}, \vec{p}_{i+1\dots n}, n) + U(\vec{q}_{1\dots i}, \vec{q}_{i+1\dots n}, n) + U(\vec{p}_{1\dots n}, \vec{q}_{1\dots n}, 2n)} \right) \\ &< 4\alpha^{z_{\vec{p}}+z_{\vec{q}}-z_{\vec{\eta}_e(\vec{p}, \vec{q})}-z_{\vec{a}'}} \end{aligned}$$

where $U(\vec{p}_{1\dots j}, \vec{q}_{i'+1\dots j'}, n) = \exp(-\sigma(m_{\vec{p}_{1\dots j}} + m_{\vec{q}_{i'+1\dots j'}} - n))$. The last inequality follows from each term in the numerator being strictly less than the entire denominator. Now let C be the exponent of the α term. Then we have¹⁶

$$\begin{aligned} C &= z_{\vec{p}} + z_{\vec{q}} - z_{\vec{\eta}_e(\vec{p}, \vec{q})} - z_{\vec{a}'} \\ &= \sum_{\vec{x} \in M} \frac{2\ell_{\vec{x}}}{1 + \exp(-\sigma(m_{\vec{x}, \vec{p}_{1\dots i}} + m_{\vec{x}, \vec{p}_{i+1\dots n}} - n))} \\ &\quad + \frac{2\ell_{\vec{x}}}{1 + \exp(-\sigma(m_{\vec{x}, \vec{q}_{1\dots i}} + m_{\vec{x}, \vec{q}_{i+1\dots n}} - n))} \\ &\quad - \frac{2\ell_{\vec{x}}}{1 + \exp(-\sigma(m_{\vec{x}, \vec{p}_{1\dots i}} + m_{\vec{x}, \vec{q}_{i+1\dots n}} - n))} \\ &\quad - \frac{2\ell_{\vec{x}}}{1 + \exp(-\sigma(m_{\vec{x}, \vec{q}_{1\dots i}} + m_{\vec{x}, \vec{p}_{i+1\dots n}} - n))}. \end{aligned}$$

Each term of the above summation is between -1 and 1 , so a worst-case upper bound is $|M|$. \square

Lemma 13. For all neighbors \vec{p} and $\vec{q} \in \Omega$,

$$\max\{\pi_\alpha(\vec{p})/\pi_\alpha(\vec{q}), \pi_\alpha(\vec{q})/\pi_\alpha(\vec{p})\} \leq \alpha^{|M|} 60Kn \ln k.$$

Proof. Since \vec{p} and \vec{q} are neighbors, they only differ in one position, so $|m_{\vec{p}} - m_{\vec{q}}| \leq 1$. Then

$$\begin{aligned} \frac{\pi_\alpha(\vec{p})}{\pi_\alpha(\vec{q})} &= \left(\frac{p(\vec{x})}{q(\vec{x})} \right) \alpha^{z_{\vec{p}}-z_{\vec{q}}} = \left(\frac{1 + \exp(\sigma n - \sigma m_{\vec{q}})}{1 + \exp(\sigma n - \sigma m_{\vec{p}})} \right) \alpha^{z_{\vec{p}}-z_{\vec{q}}} \\ &\leq \left(\frac{1 + \exp(\sigma n - \sigma(m_{\vec{p}} - 1))}{1 + \exp(\sigma n - \sigma m_{\vec{p}})} \right) \alpha^{z_{\vec{p}}-z_{\vec{q}}} \\ &= \left(\frac{1 + e^\sigma \exp(\sigma n - \sigma m_{\vec{p}})}{1 + \exp(\sigma n - \sigma m_{\vec{p}})} \right) \alpha^{z_{\vec{p}}-z_{\vec{q}}} \end{aligned}$$

¹⁶When the subscript \vec{x} is omitted from m , then m counts the number of matches with the current example. In the summations over $\vec{x} \in M$, $m_{\vec{x}}$ represents the number of matches with example \vec{x} .

$$\leq \left(\frac{1 + \exp(\sigma n - \sigma m_{\vec{p}})}{1 + \exp(\sigma n - \sigma m_{\vec{q}})} \right) e^{\sigma \alpha^{\vec{p}} - \alpha^{\vec{q}}} = e^{\sigma \alpha^{\vec{p}} - \alpha^{\vec{q}}}.$$

We now consider $z_{\vec{p}} - z_{\vec{q}}$, which equals

$$\begin{aligned} & \sum_{\vec{x} \in M} \frac{2\ell_{\vec{x}}}{1 + \exp(\sigma n - \sigma m_{\vec{x}, \vec{p}})} - \frac{2\ell_{\vec{x}}}{1 + \exp(\sigma n - \sigma m_{\vec{x}, \vec{q}})} \\ &= 2 \sum_{\vec{x} \in M} \ell_{\vec{x}} \left(\frac{\exp(\sigma n - \sigma m_{\vec{x}, \vec{q}}) - \exp(\sigma n - \sigma m_{\vec{x}, \vec{p}})}{1 + \exp(\sigma n - \sigma m_{\vec{x}, \vec{q}}) + \exp(\sigma n - \sigma m_{\vec{x}, \vec{p}}) + \exp(2\sigma n - \sigma(m_{\vec{x}, \vec{p}} + m_{\vec{x}, \vec{q}}))} \right) \\ &\leq 2 \sum_{\vec{x} \in M} \ell_{\vec{x}} \left(\frac{\exp(\sigma n - \sigma m_{\vec{x}, \vec{p}} + \sigma) - \exp(\sigma n - \sigma m_{\vec{x}, \vec{p}})}{1 + \exp(\sigma n - \sigma m_{\vec{x}, \vec{p}} + \sigma) + \exp(\sigma n - \sigma m_{\vec{x}, \vec{p}}) + \exp(2\sigma n - \sigma(m_{\vec{x}, \vec{p}} - 1 + m_{\vec{x}, \vec{p}}))} \right) \\ &= 2 \sum_{\vec{x} \in M} \ell_{\vec{x}} \left(\frac{\exp(\sigma n - \sigma m_{\vec{x}, \vec{p}})(e^{\sigma} - 1)}{1 + \exp(\sigma n - \sigma m_{\vec{x}, \vec{p}})(e^{\sigma} + 1) + e^{\sigma} \exp(2\sigma(n - m_{\vec{x}, \vec{p}}))} \right) \\ &\leq 2 \sum_{\vec{x} \in M} \frac{\ell_{\vec{x}}}{1 + \exp(\sigma(n - m_{\vec{x}, \vec{p}}))} \leq |M|, \end{aligned}$$

where the third line follows from the fact that there are only three ways to relate $m_{\vec{x}, \vec{p}}$ and $m_{\vec{x}, \vec{q}}$ for a specific \vec{x} . If $m_{\vec{x}, \vec{p}} = m_{\vec{x}, \vec{q}} + 1$, then that term of the summation equals the bound. If $m_{\vec{x}, \vec{p}} = m_{\vec{x}, \vec{q}} - 1$, then that term of the summation is negative, which is less than the bound. If $m_{\vec{x}, \vec{p}} = m_{\vec{x}, \vec{q}}$, then that term of the summation is 0, which is less than the bound.

Finally, we note that a symmetric argument can be made for $\pi_{\alpha}(\vec{q})/\pi_{\alpha}(\vec{p})$. \square

We now apply Theorem 5 to bound the mixing time of this Markov chain.

Corollary 14. *When learning generalized DNF using Winnow and a logistic $p(\vec{x})$, a simulation of \mathcal{M} that starts at any node and is of length*

$$T_i = 480kn^3\alpha_i^{1+2|M|}K \ln k(n \ln k + 2|M| \ln \alpha_i + \ln(1/\epsilon'))$$

(where $|M|$ is the number of prediction mistakes made so far) will draw samples from $\hat{\pi}_{\alpha_i}$ such that $\|\hat{\pi}_{\alpha_i} - \pi_{\alpha_i}\| \leq \epsilon'$.

Proof. Lemmas 12 and 13 bound g and h , which we substitute directly into Theorem 5. Also, note that $W(\alpha_i) \leq k^n \alpha_i^{|M|}$ and $w_{\alpha_i, \vec{p}} \geq \alpha_i^{-|M|}$, completing the proof. \square

Before many prediction mistakes have been made, our algorithm can quickly generate random samples from \mathcal{M} almost according to its stationary distribution π . Unfortunately, our worst-case mixing time bound of this chain (once $|M|$ approaches Winnow's mistake bound) is exponential in n and K . Indeed, a straightforward brute-force means to compute the sum of the weights can be done in $\Theta(k^n)$ time, whereas our bound of T_i grows with $\alpha_i^{1+2(8.88+15.54Kn \ln k)} \geq k^{31Kn \ln \alpha_i} \geq k^{12.5n}$

when $\alpha_i = \alpha = 3/2$ (a popular value for α). However, Corollary 14 is based on worst-case, adversary-based analyses. In particular, the proofs of Lemmas 12 and 13 both bound the exponents of the α terms with $|M|$ even though it is possible that they are much smaller. (For example, in Lemma 13's proof, the terms in the summation of $z_{\vec{p}} - z_{\vec{q}}$ are exponentially small when $m_{\vec{x}, \vec{p}}$ is small, which can occur frequently for terms \vec{p} with few zeroes. Also, we assumed that each term of the summation was positive, even though several could be negative.) It is open whether sub-exponential bounds can be achieved by applying a different analysis to some special cases of restricted concept classes and distributional assumptions. Further, in Section 5 we show that in practice, our algorithm performs much better than the worst-case theoretical results imply, especially considering that highly accurate estimates of the weighted sums are not needed so long as we know which side of the threshold the sum lies on.

4.1.1.3. Linear $p(\vec{x})$. Applying Theorem 7 to the linear case is not as straightforward as it was for the logistic case. As in the proof of Lemma 11, we set the entries in $\vec{\mu}$ corresponding to relevant terms to μ and the remaining entries to 0. When $\ell = 1$, in the worst case exactly one relevant term matches all n variables in \vec{x} and the remaining $K - 1$ relevant terms match 0. Then we get

$$\vec{\mu} \cdot \vec{Y} = \mu \left(1 + \frac{K-1}{n+1} \right) = \mu \left(\frac{n+K}{n+1} \right),$$

which has to be ≥ 1 . When $\ell = 0$, the worst case has all K relevant terms matching $n - 1$ variables of \vec{x} , yielding

$$\vec{\mu} \cdot \vec{Y} = \mu \left(\frac{nK}{n+1} \right) \geq \mu \left(\frac{n+K}{n+1} \right)$$

for $n, K \geq 2$. Thus it is impossible to get a $\rho > 0$ for the linear case unless we treat such worst-case examples as noise and use a non-zero \vec{E} . Following this idea, we assume that when $\ell = 1$, all relevant inputs to Winnow are at least $\gamma_1/(n+1)$ (i.e. all relevant inputs match at least $\gamma_1 - 1$ variables in \vec{x}), with of course at least one such input = 1. Further, we assume that when $\ell = 0$, all relevant inputs are at most $\gamma_0/(n+1)$. Then it is easy to show that setting

$$\mu = \frac{n+1}{n+1 + \gamma_1(K-1)}$$

and

$$\rho = \frac{n+1 + \gamma_1(K-1) - \gamma_0 K}{n+1 + \gamma_1(K-1)}$$

satisfies the conditions of Theorem 7. For example, if $\gamma_1 = 3n/4$ and $\gamma_0 = n/2$, then $\mu = (4n+4)/(n+3nK+4)$, $\rho = (n+nK+4)/(n+3nK+4) \geq 1/3$, and Theorem 7's mistake bound is

$$\begin{aligned} |M| &\leq 72 + 126 \left(\frac{4K(n+1)}{n+3Kn+4} \right) \ln \left(\frac{4N(n+1)}{n+3Kn+4} \right) + 36 \sum_{j=1}^t \vec{\mu} \cdot \vec{E}_j \\ &\leq 72 + 252 \ln N + 36 \sum_{j=1}^t \vec{\mu} \cdot \vec{E}_j \leq 72 + 252n \ln k + 36 \sum_{j=1}^t \vec{\mu} \cdot \vec{E}_j, \end{aligned} \quad (10)$$

if $n \geq 2$. Of course, this is of little value as an adversarial bound, since the third term is summed over all trials and an adversary could make each term of this summation positive. But a bound of this form might be useful under appropriate distributional assumptions.

We now bound the mixing time for the linear case.

Lemma 15. *Let $\Omega = \{0, \dots, k\}^n$. Then for all $\vec{p}, \vec{q} \in \Omega$,*

$$\pi_\alpha(\vec{p})\pi_\alpha(\vec{q}) \leq \frac{(1 + n/2)^2}{n+1} \pi_\alpha(\vec{a}') \pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q})).$$

Proof. Using the same notation introduced in the first paragraph of Lemma 12's proof, we get

$$\begin{aligned} \frac{\pi_\alpha(\vec{p})\pi_\alpha(\vec{q})}{\pi_\alpha(\vec{a}') \pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q}))} &= \frac{p(\vec{x})q(\vec{x})\alpha^{z_{\vec{p}}+z_{\vec{q}}}}{\vec{\eta}_e(\vec{p}, \vec{q})(\vec{x})\vec{a}'(\vec{x})\alpha^{z_{\vec{\eta}_e(\vec{p}, \vec{q})}+z_{\vec{a}'}}} \\ &= \left(\frac{\alpha^{z_{\vec{p}}+z_{\vec{q}}}}{\alpha^{z_{\vec{\eta}_e(\vec{p}, \vec{q})}+z_{\vec{a}'}}} \right) \left(\frac{(1 + m_{\vec{p}_{1\dots i}} + m_{\vec{p}_{i+1\dots n}})(1 + m_{\vec{q}_{1\dots i}} + m_{\vec{q}_{i+1\dots n}})}{(1 + m_{\vec{p}_{1\dots i}} + m_{\vec{q}_{i+1\dots n}})(1 + m_{\vec{q}_{1\dots i}} + m_{\vec{p}_{i+1\dots n}})} \right) \\ &\leq \alpha^{z_{\vec{p}}+z_{\vec{q}}-z_{\vec{\eta}_e(\vec{p}, \vec{q})}-z_{\vec{a}'}} \left(\frac{(1 + n/2)^2}{n+1} \right). \end{aligned}$$

The last inequality holds since the second term is maximized by setting $m_{\vec{p}_{1\dots i}} = m_{\vec{q}_{i+1\dots n}} = 0$ and $m_{\vec{q}_{1\dots i}} = m_{\vec{p}_{i+1\dots n}} = n/2$. Now let C be the first term's exponent:

$$\begin{aligned} C &= z_{\vec{p}} + z_{\vec{q}} - z_{\vec{\eta}_e(\vec{p}, \vec{q})} - z_{\vec{a}'} \\ &= \sum_{\vec{x} \in M} \ell_{\vec{x}} \left(\frac{1 + m_{\vec{x}, \vec{p}_{1\dots i}} + m_{\vec{x}, \vec{p}_{i+1\dots n}} + 1 + m_{\vec{x}, \vec{q}_{1\dots i}} + m_{\vec{x}, \vec{q}_{i+1\dots n}}}{n+1} \right. \\ &\quad \left. - \frac{1 + m_{\vec{x}, \vec{p}_{1\dots i}} + m_{\vec{x}, \vec{q}_{i+1\dots n}} + 1 + m_{\vec{x}, \vec{q}_{1\dots i}} + m_{\vec{x}, \vec{p}_{i+1\dots n}}}{n+1} \right) \\ &= 0. \quad \square \end{aligned}$$

Lemma 16. *For all neighbors \vec{p} and $\vec{q} \in \Omega$,*

$$\max\{\pi_\alpha(\vec{p})/\pi_\alpha(\vec{q}), \pi_\alpha(\vec{q})/\pi_\alpha(\vec{p})\} \leq 2\alpha^{|M|/(n+1)}.$$

Proof. Since \vec{p} and \vec{q} are neighbors, they only differ in one position, so $|m_{\vec{p}} - m_{\vec{q}}| \leq 1$. Then

$$\frac{\pi_\alpha(\vec{p})}{\pi_\alpha(\vec{q})} = \left(\frac{p(\vec{x})}{q(\vec{x})} \right) \alpha^{z_{\vec{p}}-z_{\vec{q}}} = \left(\frac{1 + m_{\vec{p}}}{1 + m_{\vec{q}}} \right) \alpha^{z_{\vec{p}}-z_{\vec{q}}} \leq \left(\frac{2 + m_{\vec{q}}}{1 + m_{\vec{q}}} \right) \alpha^{z_{\vec{p}}-z_{\vec{q}}} \leq 2\alpha^{z_{\vec{p}}-z_{\vec{q}}}.$$

We now consider $z_{\vec{p}} - z_{\vec{q}}$:

$$z_{\vec{p}} - z_{\vec{q}} = \frac{1}{n+1} \sum_{\vec{x} \in M} \ell_{\vec{x}} (1 + m_{\vec{x}, \vec{p}} - 1 - m_{\vec{x}, \vec{q}}) \leq |M|/(n+1).$$

Finally, we note that a symmetric argument can be made for $\pi_\alpha(\vec{q})/\pi_\alpha(\vec{p})$. \square

Corollary 17. *When learning generalized DNF using Winnow and a linear $p(\vec{x})$, a simulation of \mathcal{M} that starts at any node and is of length*

$$T_i = 4kn(n^2/4 + n + 1)\alpha_i^{|M|/(n+1)}(n \ln k + 2|M| \ln \alpha_i + \ln(1/\varepsilon'))$$

will draw samples from $\hat{\pi}_{\alpha_i}$ such that $\|\hat{\pi}_{\alpha_i} - \pi_{\alpha_i}\| \leq \varepsilon'$.

Proof. Lemmas 15 and 16 bound g and h , which we substitute directly into Theorem 5. Also, note that $W(\alpha_i) \leq k^n \alpha_i^{|M|}$ and $w_{\alpha_i, \vec{p}} \geq \alpha_i^{-|M|}$, completing the proof. \square

Note that if $|M| = O(n \log K)$ (either if Winnow is in the early stages of learning or if the third term of Eq. (10) is $O(n \log K)$), then the chain's mixing time is polynomial in all relevant parameters if k (the number of values each variable can take on) is a constant, yielding an FPRAS under the conditions of Corollary 9.

4.1.2. Perceptron

We now consider applying our technique to Rosenblatt's Perceptron [33] algorithm. The purpose of our analysis is primarily for contrast to the Winnow case since¹⁷ Khardon et al. [18] give a kernel function to efficiently exactly compute the weighted sums when applying Perceptron to learning DNF. But they also give an exponential lower bound on the number of mistakes that kernel perceptron makes in learning DNF: $2^{\Omega(n)}$. Thus their results do not imply an efficient DNF-learning algorithm.

We refer the reader to Section 2.1 for an overview of the Perceptron algorithm. Recall from Section 3 that if the initial weight vector is the all 0s vector, then term \vec{p} 's weight is $w_{\vec{p}} = \alpha z_{\vec{p}}$, where $z_{\vec{p}} = \sum_{\vec{x} \in M} \ell_{\vec{x}} p(\vec{x})$, M is the set of examples for which a prediction mistake is made and $\ell_{\vec{x}} \in \{-1, +1\}$ is example \vec{x} 's label.

Since our technique is only capable of estimating positive functions, we cannot allow the Perceptron's weights to be negative. Thus to each weight in "standard" Perceptron, we add a positive constant c , yielding a new weight for term \vec{p} of $w_{\vec{p}} = c + \alpha z_{\vec{p}}$, where $c = 3\alpha|M| \geq 3\alpha \max_{\vec{p} \in \Omega} \{|z_{\vec{p}}|\}$. The dot product of this new weight vector with the Perceptron inputs will then be compared to the new threshold c .

We use the same definitions and algorithm (given in Section 3) that were used for Winnow, except the "base" value of α_i is now $\alpha_0 = 0$ rather than $\alpha_1 = 1$. So our weight estimate is the same as given in Eq. (1), but the product runs from $i = 1$ to r rather than starting at $i = 2$, and we multiply it by $W(\alpha_0)$. As with Winnow, this latter quantity is easily computed exactly: For binary $p(\vec{x})$, $W(\alpha_0) = c2^n$. For logistic $p(\vec{x})$, if $k_i = k$ for all i , we get

$$W(0) = c \sum_{i=0}^n 2^i \binom{n}{i} \frac{2(k-1)^{n-i}}{1 + e^{-\sigma(i-n)}},$$

¹⁷Note, however, that other applications of Perceptron for which no kernels are available might be amenable to an MCMC-based approach to estimate the dot products.

and for the linear case, we get

$$W(0) = c \sum_{i=0}^n 2^i \binom{n}{i} \frac{(i+1)(k-1)^{n-i}}{n+1} = \frac{c(k+1)^n + 2cn(k+1)^{n-1}}{n+1}.$$

Thus all three can efficiently be computed exactly.

We now make the following argument about f_i and \hat{f}_i for all three versions of $p(\vec{x})$.

Lemma 18. *When applying binary, logistic, or linear Perceptron to learn DNF, for all i ,*

$$2/3 \leq f_i, \hat{f}_i \leq 2.$$

Proof. We focus on the actual random variables, since the estimates (the “hat” variables) have the same range. Since these variables for all three versions have $p(\vec{x})$ in the numerator and denominator, they all equal

$$f_i(\vec{p}) = \frac{c + z_{\vec{p}} \alpha_{i-1}}{c + z_{\vec{p}} \alpha_i}.$$

If $z_{\vec{p}} \geq 0$, then (since $\alpha_{i-1} < \alpha_i$) obviously $f_i(\vec{p}) \leq 1$. Also, since $c \geq 2\alpha z_{\vec{p}}$ and $\alpha \geq \alpha_i$,

$$f_i(\vec{p}) \geq \frac{c + z_{\vec{p}} \alpha_{i-1}}{c + c\alpha_i/(2\alpha)} \geq \frac{c + z_{\vec{p}} \alpha_{i-1}}{3c/2} = \frac{2}{3} \left(1 + \frac{z_{\vec{p}} \alpha_{i-1}}{c} \right) \geq 2/3.$$

If $z_{\vec{p}} < 0$, then obviously $f_i(\vec{p}) > 1$. Also,

$$f_i(\vec{p}) \leq \frac{c + z_{\vec{p}} \alpha_{i-1}}{c/2} = 2 + \frac{2z_{\vec{p}} \alpha_{i-1}}{c} \leq 2.$$

Thus $f_i(\vec{p}) \in [2/3, 2]$ for all i . \square

This leads to the following corollary.

Corollary 19. *When applying Perceptron to learn generalized DNF (with $k_i = k$ for all i for the logistic and linear cases), let the sample size $S = \lceil 390r/\varepsilon^2 \rceil$ and \mathcal{M} be simulated long enough for each sample such that the variation distance between the empirical distribution and π_{α_i} is at most $\varepsilon/(15r)$. Then for any $\delta > 0$, $\hat{W}(\alpha)$ satisfies*

$$\Pr[(1 - \varepsilon)W(\alpha) \leq \hat{W}(\alpha) \leq (1 + \varepsilon)W(\alpha)] \geq 1 - \delta.$$

To bound the number of prediction mistakes the Perceptron algorithm makes in learning DNF, we apply a result from Gentile and Warmuth [12].

Theorem 20 (Gentile and Warmuth [12]). *Let $(\vec{Y}_j, \ell_j) \in [0, 1]^N \times \{0, 1\}$ for $j = 1, \dots, t$, let $\vec{\mu}$ be an arbitrary weight vector of dimension N , and let M be the set of examples on which the Perceptron algorithm makes a prediction mistake. Then the number of mistakes made by the Perceptron*

algorithm is

$$|M| \leq \left(\frac{\|\vec{\mu}\|_2 \rho}{\gamma_{\vec{\mu}, M}} \right)^2,$$

where $\|\cdot\|_2$ is the 2-norm, $\rho \geq \|\vec{Y}\|_2$ for all $\vec{Y} \in M$, and

$$\gamma_{\vec{\mu}, M} = \frac{1}{|M|} \sum_{\vec{Y}_j \in M} \ell_j \vec{\mu} \cdot \vec{Y}_j$$

is the average margin of $\vec{\mu}$.

4.1.2.1. Binary $p(\vec{x})$. Applying Theorem 20 to the binary case is straightforward. We let $\vec{\mu}$ be 0 in all places except those corresponding to a relevant attribute, which are set to 1 (so there are K 1s in $\vec{\mu}$). In addition, we add an $(N+1)$ th position to $\vec{\mu}$, setting it to $-1/2$. This position will correspond to a 1 added to each example seen by the perceptron. Thus we get $\|\vec{\mu}\|_2 = \sqrt{K+1/4}$. Since all Perceptron inputs are from $\{0, 1\}$ for the binary case, we get $\gamma_{\vec{\mu}, M} \geq 1/2$. Further, since exactly $2^n + 1$ inputs are 1 for each example, $\rho = \sqrt{2^n + 1}$. Applying Theorem 20 yields $|M| \leq (4K+1)(2^n+1)$.

As with binary $p(\vec{x})$ with Winnow, binary $p(\vec{x})$ with Perceptron is difficult to analyze to provide non-trivial bounds on the mixing time. Thus we look at the logistic and linear cases.

4.1.2.2. Logistic $p(\vec{x})$. We begin by bounding the number of mistakes logistic Perceptron will make.

Lemma 21. *When using Equation 6 with $\sigma = \ln(60Kn \ln k)$ to specify the inputs, the number of prediction mistakes made by Perceptron when learning DNF is at most*

$$(20K+5)(2+k/(3600 \ln^2 k))^n.$$

Proof. We use the same $\vec{\mu}$ as we did for the binary case, namely 1s at the K relevant positions, $-1/2$ matching the extra 1 added to each example, and 0s elsewhere. We now bound $\gamma_{\vec{\mu}, M}$. If $\ell_j = +1$ for some trial j , then at least one of the relevant terms \vec{p} must send a 1 input to Perceptron. Thus $\ell_j \vec{\mu} \cdot \vec{Y}_j \geq 1 - 1/2 = 1/2$, where \vec{Y}_j is the vector of inputs to Perceptron (i.e. the outputs of the $p(\cdot)$ functions). If $\ell_j = -1$, then in the worst case each relevant term will be almost completely satisfied by the input example \vec{x}_j , i.e. all but one variable in each relevant term will be satisfied. If this happens, then the total contribution to $\vec{\mu} \cdot \vec{Y}_j$ that comes from the K relevant terms is $2K/(1+e^\sigma)$. Adding this to the extra $-1/2$ and multiplying by $\ell_j = -1$ yields a worst-case bound of

$$\gamma_{\vec{\mu}, M} \geq \ell_j \vec{\mu} \cdot \vec{Y}_j \geq \frac{1}{2} - \frac{2K}{1+e^\sigma} = \frac{1}{2} - \frac{2K}{1+60Kn \ln k} \geq \frac{1}{2} - \frac{1}{30 \ln 2} \geq 0.45,$$

since $n \geq 1$ and $k \geq 2$.

We now bound ρ . Recall that Eq. (8) sums the $p(\vec{x})$ values for the entire set of terms. By substituting $(p(\vec{x}))^2$ for $p(\vec{x})$ in this equation and taking the square root, we exactly get the 2-norm for any input to Perceptron:

$$\begin{aligned}\|\vec{Y}\|_2 &= \sqrt{\sum_{i=0}^n 2^i \binom{n}{i} \frac{4(k-1)^{n-i}}{(1+e^{-\sigma(i-n)})^2}} \leq \sqrt{\sum_{i=0}^n 2^i \binom{n}{i} \frac{4k^{n-i}}{e^{-2\sigma(i-n)}}} \\ &= 2 \sqrt{\sum_{i=0}^n 2^i \binom{n}{i} k^{n-i} (60Kn \ln k)^{2i-2n}} \\ &= \frac{2k^{n/2}}{(60Kn \ln k)^n} \sqrt{\sum_{i=0}^n ((7200/k)n^2 K^2 \ln^2 k)^i \binom{n}{i}} \\ &= \sqrt{4 \left(\frac{(7200/k)n^2 K^2 \ln^2 k + 1^n}{(3600/k)n^2 K^2 \ln^2 k} \right)} < 2(2 + k/(3600 \ln^2 k))^{n/2},\end{aligned}$$

since $n, K \geq 1$. Thus setting $\rho = 2(2 + k/(3600 \ln^2 k))^{n/2}$ suffices. \square

We now bound the mixing time for the Markov chain.

Lemma 22. Let $\Omega = \{0, \dots, k\}^n$. Then for all $\vec{p}, \vec{q} \in \Omega$,

$$\pi_\alpha(\vec{p})\pi_\alpha(\vec{q}) \leq 16\pi_\alpha(\vec{a}')\pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q})).$$

Proof. As in the proof of Lemma 12, let $\vec{q}_{1\dots i}$ denote (q_1, \dots, q_i) , and similarly for $\vec{p}_{1\dots i}$. Then $m_{\vec{p}} = m_{\vec{p}_{1\dots i}} + m_{\vec{p}_{i+1\dots n}}$, $m_{\vec{q}} = m_{\vec{q}_{1\dots i}} + m_{\vec{q}_{i+1\dots n}}$, $m_{\vec{\eta}_e(\vec{p}, \vec{q})} = m_{\vec{p}_{1\dots i}} + m_{\vec{q}_{i+1\dots n}}$, and $m_{\vec{a}'} = m_{\vec{q}_{1\dots i}} + m_{\vec{p}_{i+1\dots n}}$. Further, all four of these values are in $\{0, \dots, n\}$. This yields

$$\begin{aligned}\frac{\pi_\alpha(\vec{p})\pi_\alpha(\vec{q})}{\pi_\alpha(\vec{a}')\pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q}))} &= \frac{p(\vec{x})q(\vec{x})(c + \alpha z_{\vec{p}})(c + \alpha z_{\vec{q}})}{\vec{\eta}_e(\vec{p}, \vec{q})(\vec{x})\vec{a}'(\vec{x})(c + \alpha z_{\vec{\eta}_e(\vec{p}, \vec{q})})(c + \alpha z_{\vec{a}'})} \\ &= \left(\frac{(c + \alpha z_{\vec{p}})(c + \alpha z_{\vec{q}})}{(c + \alpha z_{\vec{\eta}_e(\vec{p}, \vec{q})})(c + \alpha z_{\vec{a}'})} \right) \\ &\quad \times \left(\frac{1 + U(\vec{p}_{1\dots i}, \vec{q}_{i+1\dots n}, n) + U(\vec{q}_{1\dots i}, \vec{p}_{i+1\dots n}, n) + U(\vec{p}_{1\dots n}, \vec{q}_{1\dots n}, 2n)}{1 + U(\vec{p}_{1\dots i}, \vec{p}_{i+1\dots n}, n) + U(\vec{q}_{1\dots i}, \vec{q}_{i+1\dots n}, n) + U(\vec{p}_{1\dots n}, \vec{q}_{1\dots n}, 2n)} \right) \\ &< 4 \left(\frac{(c + \alpha z_{\vec{p}})(c + \alpha z_{\vec{q}})}{(c + \alpha z_{\vec{\eta}_e(\vec{p}, \vec{q})})(c + \alpha z_{\vec{a}'})} \right),\end{aligned}$$

where $U(\vec{p}_{1\dots j}, \vec{q}_{i'+1\dots j'}, n) = \exp(-\sigma(m_{\vec{p}_{1\dots j}} + m_{\vec{q}_{i'+1\dots j'}} - n))$. The last inequality comes directly from the proof of Lemma 12. We now bound the second term. Since $p(\vec{x}) \leq 1$, each (αz) summation is upper bounded by $\alpha|M|$, so the numerator is at most $(c + \alpha|M|)^2$. Meanwhile, the denominator is at least $(c - \alpha|M|)^2$. Thus since $c = 3\alpha|M|$, the second term is at most 4. \square

Lemma 23. For all neighbors \vec{p} and $\vec{q} \in \Omega$,

$$\max\{\pi_\alpha(\vec{p})/\pi_\alpha(\vec{q}), \pi_\alpha(\vec{q})/\pi_\alpha(\vec{p})\} \leq 120Kn \ln k.$$

Proof. Since \vec{p} and \vec{q} are neighbors, they only differ in one position, so $|m_{\vec{p}} - m_{\vec{q}}| \leq 1$. Then

$$\begin{aligned} \frac{\pi_\alpha(\vec{p})}{\pi_\alpha(\vec{q})} &= \left(\frac{p(\vec{x})}{q(\vec{x})} \right) \left(\frac{c + \alpha z_{\vec{p}}}{c + \alpha z_{\vec{q}}} \right) = \left(\frac{1 + \exp(\sigma n - \sigma m_{\vec{q}})}{1 + \exp(\sigma n - \sigma m_{\vec{p}})} \right) \left(\frac{c + \alpha z_{\vec{p}}}{c + \alpha z_{\vec{q}}} \right) \\ &\leq \left(\frac{1 + \exp(\sigma n - \sigma (m_{\vec{p}} - 1))}{1 + \exp(\sigma n - \sigma m_{\vec{p}})} \right) \left(\frac{c + \alpha |M|}{c - \alpha |M|} \right) \\ &= \left(\frac{1 + e^\sigma \exp(\sigma n - \sigma m_{\vec{p}})}{1 + \exp(\sigma n - \sigma m_{\vec{p}})} \right) \left(\frac{c + \alpha |M|}{c - \alpha |M|} \right) \\ &\leq \left(\frac{1 + \exp(\sigma n - \sigma m_{\vec{p}})}{1 + \exp(\sigma n - \sigma m_{\vec{p}})} \right) 2e^\sigma = 120Kn \ln k, \end{aligned}$$

since $\sigma = \ln(60Kn \ln k)$. Finally, we note that a symmetric argument can be made for $\pi_\alpha(\vec{q})/\pi_\alpha(\vec{p})$. \square

We now apply Theorem 5 to bound the mixing time of this Markov chain.

Corollary 24. When learning generalized DNF using Perceptron and a logistic $p(\vec{x})$, a simulation of \mathcal{M} that starts at any node and is of length

$$T_i = 3840kKn^3 \ln k(n \ln k + \ln(4\alpha_i |M|) + \ln(1/\epsilon'))$$

will draw samples from $\hat{\pi}_{\alpha_i}$ such that $||\hat{\pi}_{\alpha_i} - \pi_{\alpha_i}|| \leq \epsilon'$.

Proof. Lemmas 22 and 23 bound g and h , which we substitute directly into Theorem 5. Also, note that $W(\alpha_i) \leq k^n(c + \alpha_i |M|)$ and $w_{\alpha_i, \vec{p}} \geq 1$, completing the proof. \square

4.1.2.3. Linear $p(\vec{x})$. We begin by bounding the number of mistakes Perceptron will make. However, like the linear Winnow case, worst-case (adversary) bounds are not possible since the average margin could be forced to be negative. Thus we assume that the examples are such that most of them are linearly separable and have a positive average margin $\gamma_{\vec{\mu}, M}$.

Lemma 25. When using Equation 7 to specify the inputs, and if the average margin $\gamma_{\vec{\mu}, M}$ of the sequence of examples is positive, then the number of prediction mistakes made by Perceptron when learning DNF is at most

$$\frac{5K((k+1)^n + 6n(k+1)^{n-1} + 4n(n-1)(k+1)^{n-2})}{4\gamma_{\vec{\mu}, M}^2(n+1)^2}.$$

Proof. We use the same $\vec{\mu}$ as we did for the binary and logistic cases, and hence get the same 2-norm for this vector as before. To bound ρ , recall that Eq. (9) sums the $p(\vec{x})$ values for the entire set of terms. By substituting $(p(\vec{x}))^2$ for $p(\vec{x})$ in this equation and taking the square root, we exactly get the 2-norm for any input to Perceptron:

$$\begin{aligned} \|\vec{Y}\|_2 &= \sqrt{\sum_{i=0}^n 2^i \binom{n}{i} (k-1)^{n-i} \left(\frac{i+1}{n+1}\right)^2} \\ &= \frac{\sqrt{(k+1)^n + 6n(k+1)^{n-1} + 4n(n-1)(k+1)^{n-2}}}{n+1}, \end{aligned}$$

which completes the proof. \square

Now we bound the mixing time.

Lemma 26. Let $\Omega = \{0, \dots, k\}^n$. Then for all $\vec{p}, \vec{q} \in \Omega$,

$$\pi_\alpha(\vec{p}) \pi_\alpha(\vec{q}) \leq \frac{4(1+n/2)^2}{n+1} \pi_\alpha(\vec{a}') \pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q})).$$

Proof. Using the same notation introduced in the first paragraph of Lemma 22's proof, we get

$$\begin{aligned} \frac{\pi_\alpha(\vec{p}) \pi_\alpha(\vec{q})}{\pi_\alpha(\vec{a}') \pi_\alpha(\vec{\eta}_e(\vec{p}, \vec{q}))} &= \left(\frac{p(\vec{x}) q(\vec{x})}{\vec{\eta}_e(\vec{p}, \vec{q})(\vec{x}) \vec{a}'(\vec{x})} \right) \left(\frac{(c + \alpha z_{\vec{p}})(c + \alpha z_{\vec{q}})}{(c + \alpha z_{\vec{\eta}_e(\vec{p}, \vec{q})})(c + \alpha z_{\vec{a}'})} \right) \\ &\leq \left(\frac{(1+n/2)^2}{n+1} \right) \cdot 4, \end{aligned}$$

using results from the proofs of Lemmas 15 and 22. \square

Lemma 27. For all neighbors \vec{p} and $\vec{q} \in \Omega$,

$$\max\{\pi_\alpha(\vec{p})/\pi_\alpha(\vec{q}), \pi_\alpha(\vec{q})/\pi_\alpha(\vec{p})\} \leq 4.$$

Proof. Since \vec{p} and \vec{q} are neighbors, they only differ in one position, so $|m_{\vec{p}} - m_{\vec{q}}| \leq 1$. Then

$$\frac{\pi_\alpha(\vec{p})}{\pi_\alpha(\vec{q})} = \left(\frac{p(\vec{x})}{q(\vec{x})} \right) \left(\frac{c + \alpha z_{\vec{p}}}{c + \alpha z_{\vec{q}}} \right) = \left(\frac{1 + m_{\vec{p}}}{1 + m_{\vec{q}}} \right) \left(\frac{c + \alpha z_{\vec{p}}}{c + \alpha z_{\vec{q}}} \right) \leq \left(\frac{2 + m_{\vec{q}}}{1 + m_{\vec{q}}} \right) 2 \leq 4.$$

Finally, we note that a symmetric argument can be made for $\pi_\alpha(\vec{q})/\pi_\alpha(\vec{p})$. \square

Corollary 28. *When learning generalized DNF using Perceptron and a linear $p(\vec{x})$, a simulation of \mathcal{M} that starts at any node and is of length*

$$T_i = \left(\frac{32kn^2(1+n/2)^2}{(n+1)} \right) (n \ln k + \ln(4\alpha_i |M|) + \ln(1/\epsilon'))$$

will draw samples from $\hat{\pi}_{\alpha_i}$ such that $\|\hat{\pi}_{\alpha_i} - \pi_{\alpha_i}\| \leq \epsilon'$.

Proof. Lemmas 26 and 27 bound g and h , which we substitute directly into Theorem 5. Also, note that $W(\alpha_i) \leq k^n(c + \alpha_i |M|)$ and $w_{\alpha_i, \vec{p}} \geq 1$, completing the proof. \square

4.2. Pruning ensembles of classifiers

We now apply our methods to pruning an ensemble, produced by e.g. AdaBoost [35]. AdaBoost's output is a set of functions $h_i: \mathcal{X} \rightarrow \mathfrak{R}$, where $i \in \{1, \dots, n\}$ and \mathcal{X} is the instance space. Each h_i is trained on a different distribution over the training examples and is associated with a parameter $\beta_i \in \mathfrak{R}$ that weights its predictions. Given an instance $\vec{x} \in \mathcal{X}$, the ensemble's prediction is $H(\vec{x}) = \text{sign}(\sum_{i=1}^n \beta_i h_i(\vec{x}))$. Thus $\text{sign}(h_i(\vec{x}))$ is h_i 's prediction on \vec{x} , $|h_i(\vec{x})|$ is its confidence in its prediction, and β_i weights AdaBoost's confidence in h_i . It has been shown that if each h_i has error less than $1/2$ on its distribution, then the error on the training set and the generalization error of $H(\cdot)$ can be bounded. Strong bounds on $H(\cdot)$'s generalization error can also be shown even if the boosting algorithm is run past the point where $H(\cdot)$'s error is zero [34]. However, overfitting can still occur [26], i.e. sometimes better generalization can be achieved if some of the h_i 's are discarded. So our goal is to find a weighted combination of all possible prunings that performs not much worse in terms of generalization error than the best single pruning.

To predict nearly as well as the best pruning, we place every possible pruning in a pool (so $N = 2^n$) and run WM. We start by computing W^+ and W^- , which are, respectively, the sums of the weights of the experts predicting a positive and a negative label on example \vec{x} . Then WM predicts +1 if $W^+ > W^-$ and -1 otherwise. Whenever WM makes a prediction mistake, it reduces the weights of all experts that predicted incorrectly by dividing them by α (see Section 2.1).

As in Section 4.1, using a binary $p(\vec{x})$ makes bounding the mixing time difficult, except in a trivial sense. Thus we use a linear $p(\vec{x})$, which allows us to also incorporate each pruning's confidence in its prediction, and to use that confidence when updating the weights. Given an example $\vec{x} \in \mathcal{X}$, we compute $h_i(\vec{x})$ for all $i \in \{1, \dots, n\}$. We then use our MCMC procedure to compute \hat{W}^+ , an estimate of $W^+ = \sum_{\vec{p} \in \Omega^+} p(\vec{x}) w_{\vec{p}}$, where $p(\vec{x}) = \sum_{i=1}^n p_i \beta_i h_i(\vec{x})$, $\Omega^+ = \{\vec{p} \in \{0, 1\}^n: \sum_{i=1}^n p_i \beta_i h_i(\vec{x}) \geq 0\}$, $w_{\vec{p}} = \alpha^{z_{\vec{p}}}$, $z_{\vec{p}} = \sum_{\vec{x} \in M} \ell_{\vec{x}} p(\vec{x})$, and M is the set of examples for which a prediction mistake was made. A similar procedure is used to compute \hat{W}^- . Then WM predicts +1 if $\hat{W}^+ > \hat{W}^-$ and -1 otherwise.

Define the Markov chain \mathcal{M} with state space Ω^+ (similarly, Ω^-) and that makes transitions according to the description in Section 3. The chain corresponds to a random walk on the Boolean hypercube truncated by a hyperplane. It is easy to show that all pairs of states in Ω^+

(similarly, Ω^-) can communicate. To move from node $\vec{p} \in \Omega^+$ to $\vec{q} \in \Omega^+$, first add to \vec{p} all bits i in \vec{q} and not in \vec{p} that correspond to positions where $\beta_i h_i(x) \geq 0$. Then delete from \vec{p} all bits i in \vec{p} and not in \vec{q} that correspond to positions where $\beta_i h_i(x) < 0$. Then delete the unnecessary “positive bits” and add the necessary “negative bits”. It is easy to see that all states between \vec{p} and \vec{q} are in Ω^+ . Thus \mathcal{M} is irreducible and hence ergodic by Lemma 1.

As before, we let B be an upper bound¹⁸ on the sum of all updates made on any pruning. Then it is straightforward to adapt Lemma 8’s proof to bound $f, \hat{f} \in [1/e, e]$ for WM when applying Section 3’s procedure to estimate $W(\alpha)$. But when applying Eq. (1), we must determine $W(\alpha_1) = W(1) = |\Omega^+|$. This is equivalent to counting the number of solutions to a 0–1 knapsack problem, which is #P-complete. Thus in order to complete our computation of \hat{W} , we must also estimate $|\Omega^+|$. We do this by mapping the problem to the knapsack problem which is summarized in Section 2.4 and shown to have an FPRAS by Morris and Sinclair [28]. If we let the “weight” of item i in our problem be $w_i = \beta_i h_i(\vec{x})$ for an example \vec{x} , then the only difference between the two problems is that the weights in the $|\Omega^+|$ estimation problem may be negative. We now argue that they are still equivalent and thus we can directly apply the results of Morris and Sinclair. Given a vector $\vec{p} \in \{0, 1\}^n$ and a weight vector \vec{w} , let $p'_i = 1 - p_i$ if $w_i < 0$ and $p'_i = p_i$ otherwise. Also, let $w'_i = |w_i|$ and $b' = \sum_{w_i < 0} |w_i|$. It is easy to argue that $\sum_{i=1}^n w_i p_i \geq 0$ (which is the definition¹⁹ of Ω^+) if and only if $\sum_{i=1}^n w'_i p'_i \geq b'$ (which is an instance of the knapsack problem). If we let $s^+ = \sum_{w_i > 0, p_i=1} w_i$, $s^- = \sum_{w_i < 0, p_i=1} w_i$, and $s'^- = \sum_{w_i < 0, p'_i=0} w_i$, then $b' = s'^- - s^-$. This is exactly what is added to both sides of the inequality $\sum_{i=1}^n w_i p_i \geq 0$ to get $\sum_{i=1}^n w'_i p'_i \geq b'$. Thus we can efficiently estimate $|\Omega^+|$ to within a factor of ε , allowing us to apply Corollary 3.

Corollary 29. *When applying WM to learn a weighted combination of ensemble prunings, let the sample size $S = \lceil 130 re^2/\varepsilon^2 \rceil$, $|\Omega^+|$ be estimated to within $\varepsilon/2$ of its true value with probability $\geq 3/4$ via the procedure outlined above, and \mathcal{M} be simulated long enough for each sample such that the variation distance between the empirical distribution and π_{α_i} is at most $\varepsilon/(10e^2r)$. Then for any $\delta > 0$, $\hat{W}^+(\alpha)$ satisfies*

$$\Pr[(1 - \varepsilon)W^+(\alpha) \leq \hat{W}^+(\alpha) \leq (1 + \varepsilon)W^+(\alpha)] \geq 1 - \delta,$$

and the same result applies to $\hat{W}^-(\alpha)$ if $|\Omega^-|$ is well approximated.

Note that if $W^+/W^- \notin [\frac{1-\varepsilon}{1+\varepsilon}, \frac{1+\varepsilon}{1-\varepsilon}]$ for all trials, then our estimates of W^+ and W^- are (with probability at least $1 - \delta'_t$ for trial t) sufficiently accurate to correctly determine whether or not $W^+ > W^-$. Setting $\delta'_t = \delta/2^t$ yields a total probability of failure of at most²⁰ $\sum_{t=1}^{\infty} \delta/2^t = \delta$. Thus

¹⁸In contrast to Section 4.1, where B is by definition upper bounded by Winnow’s or Perceptron’s mistake bound, for this application B could be arbitrarily large since it depends on the predictions of arbitrary hypotheses. Thus for the rest of this section we implicitly assume that B is polynomial in all relevant parameters, i.e. that it is expressed in unary.

¹⁹By negating all w_i , we can use the same arguments to estimate $|\Omega^-|$.

²⁰Recall from the proof of Theorem 2 that only $O(\log 1/\delta')$ runs of the estimation procedure are needed to reduce the probability of failure to δ' .

under these conditions, our version of WM runs identically to the brute-force version, and we can apply WM's mistake bounds. This yields the following corollary.

Corollary 30. *Using the assumptions of Corollary 29, if $W^+/W^- \notin [\frac{1+\epsilon}{1-\epsilon}, \frac{1-\epsilon}{1+\epsilon}]$ for all t , then with probability at least $1 - \delta$, the number of prediction mistakes made by this algorithm on any sequence of examples is $O(v + n)$, where n is the number of hypotheses in the ensemble and v is the number of mistakes made by the best pruning.*

Now we bound the mixing time. Bounding g is easy, since when viewed as multisets, $\vec{p} \cup \vec{q} = \vec{a}' \cup \vec{\eta}_e(\vec{p}, \vec{q})$, which implies $z_{\vec{p}} + z_{\vec{q}} = z_{\vec{a}'} + z_{\vec{\eta}_e(\vec{p}, \vec{q})}$ and $\pi_{\alpha}^+(\vec{p})\pi_{\alpha}^+(\vec{q}) = \pi_{\alpha}^+(\vec{a}')\pi_{\alpha}^+(\vec{\eta}_e(\vec{p}, \vec{q}))$. Thus $g = 1$. Further, if \vec{p} and \vec{q} are neighbors that differ in bit i , then

$$\begin{aligned} z_{\vec{p}} - z_{\vec{q}} &= \sum_{\vec{x} \in M} \left(\sum_{j:p_j=1} \ell_{\vec{x}} \beta_j h_j(\vec{x}) - \sum_{j:q_j=1} \ell_{\vec{x}} \beta_j h_j(\vec{x}) \right) \\ &\leq \beta_i \sum_{\vec{x} \in M} \ell_{\vec{x}} h_i(\vec{x}), \end{aligned}$$

which implies that for any neighbors \vec{p} and \vec{q} , $\max\{\pi_{\alpha}(\vec{p})/\pi_{\alpha}(\vec{q}), \pi_{\alpha}(\vec{q})/\pi_{\alpha}(\vec{p})\} \leq \alpha^{B_{\max}}$ where $B_{\max} = \max_j \{\beta_j \sum_{\vec{x} \in M} \ell_{\vec{x}} h_j(\vec{x})\}$.

Corollary 31. *When learning a weighted combination of ensemble prunings with Weighted Majority, if $\Omega = \{0, 1\}^n$, then a simulation of \mathcal{M} that starts at any node and is of length*

$$T_i = 2n^2 \alpha_i^{B_{\max}} (n \ln 2 + n(B_{\max} - B_{\min}) \ln \alpha_i + \ln 1/\epsilon')$$

will draw samples from $\hat{\pi}_{\alpha_i}$ such that $\|\hat{\pi}_{\alpha_i} - \pi_{\alpha_i}\| \leq \epsilon'$.

Proof. We substitute our bounds of g and h directly into Theorem 5 and note that $W(\alpha_i) \leq 2^n \alpha_i^{nB_{\max}}$ and $w_{\alpha_i, \vec{p}} \geq \alpha_i^{nB_{\min}}$, completing the proof. \square

The above mixing time bound is only polynomial if B_{\max} is logarithmic in all relevant parameters, which is unlikely. However, our analysis is handicapped by worst-case assumptions, as with Corollary 14. While it is unlikely that an efficient bound on the mixing time exists for general ensembles with arbitrary classifiers and arbitrary distributions over examples, it is open whether restricted cases could have better bounds. Further, in Section 5 we show that in practice, our algorithm performs much better than the worst-case theoretical results imply, especially considering that highly accurate estimates of the weighted sums are not needed so long as we know whether or not $W^+ > W^-$.

Note that in Corollary 31 we assume that $\Omega = \{0, 1\}^n$, i.e. that all prunings classify \vec{x} as positive (or negative), and the chain is an untruncated hypercube. This is because without such an assumption we cannot guarantee that a canonical path between two nodes does not leave Ω at any time. A new approach, employing *balanced, almost-uniform permutations* has recently been

pioneered by Morris and Sinclair [28] and applied to the truncated Boolean hypercube problem when the chain's stationary distribution is uniform (i.e. for counting the number of solutions to the knapsack problem). It is reasonable that their technique could be generalized to the case of a non-uniform distribution, which would allow us to consider truncated hypercubes for WM and for other algorithms.

4.3. Discussion

In examining the results and proofs related to using Winnow on DNF, we see some interesting differences. Recall from Theorem 5 that there are two functions that must be bounded in order to bound a chain's mixing time: g , which bounds the ratio of $\pi(\vec{p})\pi(\vec{q})$ to $\pi(\vec{d}')\pi(\vec{\eta}_e(\vec{p}, \vec{q}))$, and h , which bounds the ratio of weights of neighboring nodes in the chain. Linear $p(\vec{x})$ allows us to bound g with a polynomial since the z 's in the exponent of α all cancel out due to the linear nature of the weight updates. However, logistic $p(\vec{x})$ (and probably binary $p(\vec{x})$) in the worst case get charged a multiplicative factor of α for each update, yielding an upper bound of g that is exponential in $|M|$. In contrast, when bounding h in an adversarial setting, both linear and logistic $p(\vec{x})$ allow the difference in the z 's to grow with $|M|$. Further, in the worst case we cannot bound $|M|$ for linear, but we can for logistic and binary.

A more fundamental difference arises when comparing multiplicative weight update algorithms (Winnow and WM) to the additive weight update algorithm (Perceptron). The additive weight updates prevent \mathcal{M} 's stationary distribution from deviating very far from uniform (when compared to the MWU algorithms). Thus for Perceptron, \mathcal{M} mixes rapidly, even if $|M|$ is exponentially large. However, from a learning-theoretic standpoint, Perceptron is not a good choice for learning DNF since Lemma 21's bound on the number of mistakes (updates) Perceptron will make in the worst case is exponential, which is corroborated by Kharton et al.'s lower bound [18].

A natural extension of the results of Section 4.2 is to generalize the results to multiclass predictions. This comes for free if the boosting algorithm used is AdaBoost.MH from Schapire and Singer [35], which reduces the multiclass boosting problem to a set of binary ones, allowing our results to fit into this framework. Alternatively, AdaBoost.M1 from Freund and Schapire [11] more directly addresses the multiclass problem by having each hypothesis predict its confidence that \vec{x} belongs to class j . Then the ensemble's prediction is the class that maximizes these confidence-rated predictions. That is, each class is tested individually and the one that scores the highest is the predicted class for \vec{x} . To adapt our framework to this, rather than simply estimating W^+ and W^- , we estimate W^j for each class j and then predict the class with the maximum. Since each W^j estimate uses a separate Boolean hypercube truncated by a single hyperplane (the one that separates prunings that predict class j from those predicting another class), we can bound the mixing time using the same machinery developed in Section 4.2, assuming a version of Theorem 5 exists for truncated cubes.

Since one of the goals of pruning an ensemble of classifiers is to reduce its size, one may adopt one of several heuristics, such as choosing the pruning that has highest weight in WM, the highest ratio of weight to size, or the highest product of weight and diversity, where diversity is measured by e.g. Kullback–Leibler divergence (see [26]). Let $f(\vec{p})$ be the function that one wants to maximize. Then the goal is to find the $\vec{p} \in \{0, 1\}^n$ that approximately maximizes f . To do this one

can define a new Markov chain \mathcal{M}' whose transition probabilities are the same as for \mathcal{M} in Section 3 except that Step 3 is irrelevant (since there is no training example \vec{x}) and in Step 4, change the transition probability to $\min\{1, \rho^{f(\vec{p}') - f(\vec{p})}\}$. The parameter ρ governs the shape of the stationary distribution: $\rho = 1$ implies a uniform distribution over all prunings, while a large value of ρ yields a distribution that peaks at prunings with large $f(\vec{p})$. (This is a special case of simulated annealing [19] where the temperature is held constant.) Lemma 1 obviously holds for \mathcal{M}' , but it is an open problem to bound how far from optimal its solution will be. Of course, other combinatorial optimization methods such as genetic algorithms can also be applied here.

Similarly, one issue with our DNF algorithm is that after training, we still require the training examples and running \mathcal{M} to evaluate the hypothesis on a new example. In lieu of this, one can, after training, search (using a modified chain or a GA as described above) for the terms with the largest weights in Winnow. The result is a set of rules, and the prediction on a new example can be a thresholded sum of weights of satisfied rules, using the same threshold θ . The only issue then is to determine how many terms to select. Since each example satisfies exactly 2^n terms, for an example to be classified as positive, the average weight of its satisfied terms must be at least $\theta/2^n$. Thus one heuristic is to choose as many terms as possible with weight at least $\theta/2^n$, stopping when we find “too many” (as specified by a parameter) terms with weight less than $\theta/2^n$. Using this pruned set of rules, no additional false positives will occur, and in fact the number might be reduced. The only concern is causing extra false negatives.

5. Empirical results

The primary purpose of our experiments is to assess how well our algorithms will work in practice, especially when compared to our worst-case bounds. A more thorough empirical analysis (as well as some heuristic optimizations) of our technique is given by Tao and Scott [41].

The goal of our algorithms is to use the weighted sum approximations to accurately simulate Winnow and WM, since an accurate simulation is required for us to apply Winnow’s and WM’s error bounds.²¹ We measure accuracy of simulation in several ways: (1) comparing weighted sum estimates to their true values (computed by brute-force implementations); (2) counting the number of times our algorithm predicts differently from brute-force (e.g. for Winnow, the fraction of weighted sum estimates that are on the opposite side of the threshold as the true weighted sum); and (3) measuring prediction error. The simulated data we used in our experiments should make learning straightforward for brute-force, so low prediction error should correlate (to some extent) to simulation accuracy. This performance measure is especially useful on problems that are too large for brute-force to handle.

5.1. Learning DNF formulas

In our DNF experiments, we defined the instance space to be $\mathcal{X} = \{1, 2\}^n$ and the set of terms to be $\mathcal{P} = \{0, 1, 2\}^n$, i.e. $k_i = 2$ for all i . Recall that a term $\vec{p} = (p_1, \dots, p_n) \in \mathcal{P}$ is satisfied by

²¹ Conceivably, in some cases our algorithm may accidentally make fewer prediction mistakes when deviating from brute-force implementations, but in the absence of a formal theory to characterize this, it is safer to assume that such behavior is anomalous. Thus our primary goal is to measure how well we simulate brute-force.

example $\vec{x} = (x_1, \dots, x_n) \in \mathcal{X}$ if and only if $\forall p_i > 0, p_i = x_i$. So $p_i = 0 \Rightarrow x_i$ is irrelevant for term \vec{p} and $p_i > 0 \Rightarrow x_i$ must equal p_i for \vec{p} to be satisfied. Even though we do not have an analysis of its mixing time, we used binary $p(\vec{x})$ in our experiments.²² Even though using binary $p(\vec{x})$ we could define for each new example $\Omega = \{0, 1\}^n$ (i.e. an untruncated Boolean hypercube; see Section 4.1.1), we chose to use as Ω a truncated (with a single hyperplane) version of \mathcal{P} . We did this for two reasons. First, doing so allows us to evaluate the performance of our MCMC approach when the hypercube is truncated, which is more generally applicable and also currently lacking in theoretical results. Second, this experimental approach gives us a better idea of how quickly the time complexity of a brute-force implementation grows as a function of n . Comparing this with the time of our MCMC experiments tells us the minimum value of n for which our approach is faster.

We generated random (from a uniform distribution) $K = 5$ -term DNF formulas, using $n \in \{10, 15, 20\}$. So the total number of Winnow inputs was $3^{10} = 59049$, $3^{15} = 1.43 \times 10^7$, and $3^{20} = 3.49 \times 10^9$. For each value of n there were nine training/testing set combinations, each with 50 training examples and 50 testing examples. Examples were generated uniformly at random.

Table 1 gives averaged²³ results for $n = 10$, indexed by S and T (“BF” means brute-force). “GUESS” is the average error of the estimates ($= |\text{guess} - \text{actual}| / \text{actual}$). “LOW” is the fraction of guesses that were $< \theta$ when the actual value was $> \theta$, and “HIGH” is symmetric. These are the only times our algorithm deviates from brute-force. “PM” is the number of prediction mistakes made by Winnow on the training set while learning (in all experiments, Winnow repeatedly made passes over the training set until it correctly classified all training examples). This gives an evaluation of each algorithm in an on-line setting. After training was complete, we also evaluated the hypotheses on their respective test sets, i.e. in a batch learning setting. “GE” is the generalization error on the test set. Finally, “ S_{theo} ” and “ T_{theo} ” are S and T from Corollaries 9 and 14 that guarantee an error of GUESS given the values of r in our simulations using $\alpha = 3/2$. These latter two columns show how pessimistic the worst-case bounds are in contrast to what works in practice. In general, the results in the columns varied little across the different data sets: the standard deviations of the values were typically small when compared to the means.

Both GUESS and HIGH are very sensitive to T but not as sensitive to S . LOW was negligible due to the distribution of weights as training progressed: the term $\vec{p}_e = \vec{0}$ (satisfied by all examples) had high weights. Since all computations started at $\vec{0}$ and the Markov chain \mathcal{M} seeks out nodes with high weights, the estimates tended to be too high rather than too low. But this is less significant as S and T increase. For $S = 100$ and $T = 300$, training and testing with \mathcal{M} was slower than brute-force by a factor of over 108. The average value of r used was 20.79 (range was 19–26).

Since the run time of our algorithm varies linearly in r , we ran some experiments where we fixed r rather than letting it be set as in Section 4.1. We set $S = 100$, $T = 300$ and $r \in \{5, 10, 15, 20\}$. The results are in Table 2. This indicates that for the given parameter values, r can be reduced below that which is stipulated in Section 3.

²²When we compare actual mixing times to theoretical bounds, we will compare to theoretical bounds for logistic $p(\vec{x})$, which is a reasonable approximation to the binary case.

²³The number of weight estimations made per row in the table varied due to a varying number of training rounds, but typically was around 3000.

Table 1

Results for $n = 10$ and r chosen as in Section 3

S	T	GUESS	LOW	HIGH	PM	GE	S_{theo}	T_{theo}
100	100	0.4713	0.0000	0.1674	35.67	0.0600	2.23×10^5	1.772×10^{104}
100	200	0.1252	0.0017	0.0350	35.67	0.0533	3.16×10^6	1.777×10^{104}
100	300	0.0634	0.0041	0.0172	37.89	0.0711	1.23×10^7	1.780×10^{104}
100	500	0.0484	0.0091	0.0078	40.11	0.0844	2.11×10^7	1.781×10^{104}
500	100	0.4826	0.0000	0.1594	34.67	0.1000	2.13×10^5	1.772×10^{104}
500	200	0.1174	0.0000	0.0314	33.83	0.0600	3.60×10^6	1.778×10^{104}
500	300	0.0441	0.0043	0.0145	34.22	0.0867	2.55×10^7	1.781×10^{104}
500	500	0.0232	0.0034	0.0064	37.88	0.0800	9.16×10^7	1.784×10^{104}
BF					36.56	0.0730		

Table 2

Results for $n = 10$, $S = 100$, and $T = 300$

r	GUESS	LOW	HIGH	PM	GE
5	0.1279	0.0119	0.0203	40.67	0.0844
10	0.0837	0.0095	0.0189	38.33	0.0867
15	0.0711	0.0058	0.0159	37.78	0.0800
20	0.0638	0.0042	0.0127	36.22	0.0889
BF				36.56	0.0730

Table 3

Results for $n = 15$ and r chosen as in Section 3

S	T	GUESS	LOW	HIGH	PM	GE	S_{theo}	T_{theo}
500	1500	0.0368	0.0028	0.0099	60.22	0.0700	5.01×10^7	4.112×10^{151}
500	1800	0.0333	0.0040	0.0049	60.75	0.0675	6.12×10^7	4.112×10^{151}
500	2000	0.0296	0.0035	0.0023	57.00	0.0675	7.68×10^7	4.113×10^{151}
1000	1500	0.0388	0.0015	0.0042	56.25	0.0650	4.51×10^7	4.111×10^{151}
1000	1800	0.0253	0.0006	0.0038	59.00	0.0775	1.06×10^8	4.114×10^{151}
1000	2000	0.0207	0.0025	0.0020	49.00	0.0800	1.58×10^8	4.115×10^{151}
BF					60.22	0.0800		

Results for $n = 15$ appear in Table 3. The trends for $n = 15$ are similar to those for $n = 10$. Brute-force is faster than \mathcal{M} at $S = 500$ and $T = 1500$, but only by a factor of 16. The average value of r used was 31.52 (range was 19–40). As with $n = 10$, r can be reduced to speed up the algorithm, but at a cost of increasing the errors of the predictions (e.g. see Table 4). We ran the same experiments with a training set of size 100 rather than 50 (the test set was still of size 50), summarized in Table 5. As expected, error on the guesses changes little, but GE is decreased.

For $n = 20$, no exact (brute-force) sums were computed since there are over 3 billion inputs. So we only examined the prediction error of our algorithm. With $S = 1000$, $T = 2000$, r set as in

Table 4

Results for $n = 15$, $S = 500$, and $T = 1500$, and a training set of size 50

r	GUESS	LOW	HIGH	PM	GE
10	0.0572	0.0049	0.0132	59.22	0.1075
20	0.0444	0.0033	0.0063	61.22	0.0756
30	0.0407	0.0022	0.0047	62.00	0.0822
BF				60.22	0.0800

Table 5

Results for $n = 15$, $S = 500$, and $T = 1500$, and a training set of size 100

r	GUESS	LOW	HIGH	PM	GE
10	0.0577	0.0046	0.0478	78.00	0.0511
20	0.0456	0.0032	0.0073	78.33	0.0733
30	0.0405	0.0044	0.0081	74.44	0.0689
BF				80.11	0.0356

Section 3, and a training set of size 100, the average number of prediction mistakes was 91.75 and the average GE was 0.11. The average value of r used was 55 (range was 26–78), and the run time for \mathcal{M} was over 270 times faster than brute-force (brute-force was run on a small number of examples to estimate its time complexity for $n = 20$). Thus for this case our algorithm provides a significant speed advantage. When running our algorithm with a fixed value of $r = 30$ (reducing time per example by almost a factor of 2), GE increases to 0.1833.

In summary, even though our experiments are for small values of n , they indicate that relatively small values of S , T , and r are sufficient to minimize our algorithm's deviations from brute-force Winnow. In addition, our algorithm becomes significantly faster than that of brute-force somewhere between $n = 15$ and $n = 20$, which is small for a machine learning problem. However, our implementation is still extremely slow, taking several days or longer to finish training when $n = 20$ (evaluating the learned hypothesis is also slow). Thus we are actively working on optimizations to speed up learning and evaluation (see Section 6).

5.2. Ensemble pruning

For the Weighted Majority experiments, we used AdaBoost over decision shrubs (depth-2 decision trees) generated by C4.5 [31] to learn hypotheses for an artificial two-dimensional data set (Fig. 1). The target concept is a circle of radius 10 and the examples are distributed around its circumference, each point's distance from the circle normally distributed with zero mean and unit variance. By concentrating examples around the circular boundary and limiting each decision tree's depth, we required ensembles of multiple trees to achieve low classification error on the data. We created an ensemble of 10 classifiers and simulated WM with²⁴ $S \in \{50, 75, 100\}$ and

²⁴The estimation of $|\Omega^+|$ required an order of magnitude larger values of S and T than did the estimation of the ratios to get sufficiently low error rates.

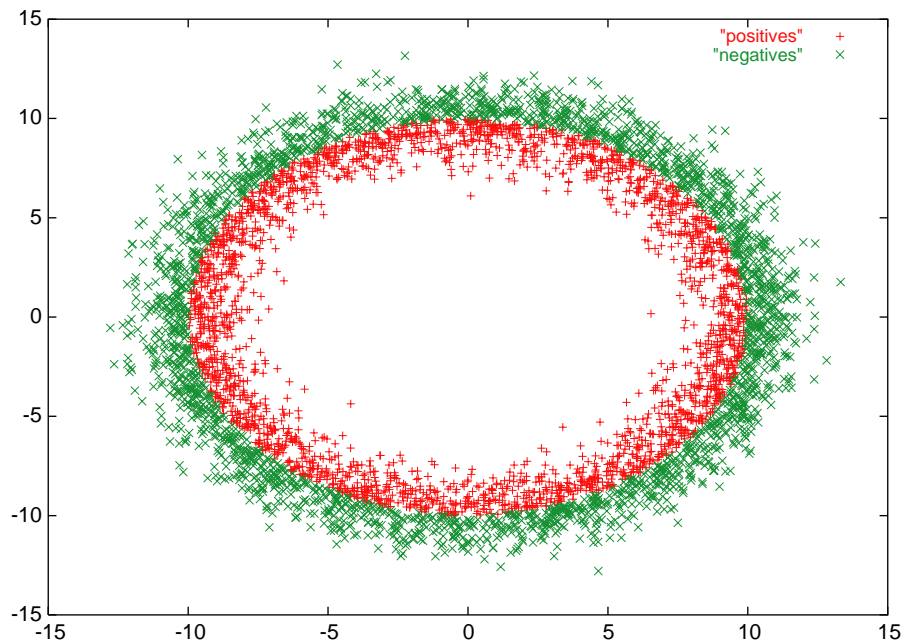


Fig. 1. The circle data set.

Table 6

Results for $n = 10$ and r chosen as in Section 3

S	T	$ \Omega^+ $	X_i	$\hat{W}^+(\alpha)$	DEPARTURE
50	500	0.0423	0.00050	0.0071	0.0000
50	750	0.0332	0.00069	0.0061	0.0000
50	1000	0.0419	0.00068	0.0070	0.0000
75	500	0.0223	0.00067	0.0050	0.0000
75	750	0.0197	0.00047	0.0047	0.0000
75	1000	0.0276	0.00058	0.0055	0.0000
100	500	0.0185	0.00040	0.0047	0.0000
100	750	0.0215	0.00055	0.0050	0.0000
100	1000	0.0288	0.00044	0.0056	0.0000

$T \in \{500, 750, 1000\}$ on the set of 2^{10} prunings and compared the values computed for Eq. (1) to the true values from brute-force WM. The results are in Table 6: “ $|\Omega^+|$ ” denotes the error of our estimates of $|\Omega^+|$, “ X_i ” denotes the error of our estimates of the ratios $W^+(\alpha_i)/W^+(\alpha_{i-1})$, and “ $\hat{W}^+(\alpha)$ ” denotes the error of our estimates of $W^+(\alpha)$. Finally, “DEPARTURE” indicates our algorithm’s departure from brute-force WM, i.e. in these experiments our algorithm perfectly

emulated brute-force. Finally, we note that other results show that for $n = 30$, $S = 200$, and $T = 2000$, our algorithm takes about 4.5 h/example to run, while brute-force takes about 2.8 h/example. Thus we expect our algorithm to run faster than brute-force at about $n = 31$ or $n = 32$.

6. Conclusions and future work

MWU algorithms are particularly useful when the number of inputs is very large, since their mistake bounds are logarithmic in the total number of inputs. However, only in specific cases is it known how to compute the weighted sum of inputs efficiently to exploit this attribute efficiency. We presented a general, widely applicable method based on Markov chain Monte Carlo to estimate these weighted sums, along with theoretical and empirical analyses of these methods as applied to learning DNF formulas and pruning ensembles of classifiers. Our theoretical results do not yield efficient algorithms for these problems, but they do provide machinery for potentially conducting average-case analyses of these algorithms on e.g. restricted classes of DNF and/or on restricted distributions. Further, as a heuristic, our methods show promise: in experimental results on simulated data, our algorithms perform much better than the worst-case theoretical results imply, especially considering that highly accurate estimates of the weighted sums are not needed so long as we know which side of the threshold the sum lies on. Also, our approach can very easily be implemented on a parallel or clustered architecture since all samples are drawn independently of each other, and each term in Eq. (1) can be computed independently of the other terms.

Recent work by Tao and Scott [41] includes a thorough empirical analysis of this method to speed it up further. They conducted tests using data sets from the UCI Repository [1] and included experimenting with other sampling methods besides the Metropolis sampler [27] of Section 3. They also utilized methods that stop sampling early when it is known what side of the threshold the weighted sum will fall.

It is open whether better mixing time bounds are possible for special cases of the DNF problem of Section 4.1. For example, if one considers learning parity functions (or other classes of functions) under the uniform distribution, can the bounds of Lemmas 12 and 13 be tightened to sub-exponential? Alternatively, are there special cases where Winnow with linear $p(\vec{x})$ (Section 4.1.1.3) has on average a sufficiently small mistake bound to make the mixing time polynomial?

It should be possible to apply our results to other problems for which an MWU algorithm is applicable with an exponential number of inputs. The key is to map the set of inputs to a hypercube (perhaps truncated), and use that space as the set of states for the Markov chain. One such application would be to accelerate Winnow-based algorithms [13,37,42] for a learning model that generalizes the conventional multiple-instance learning model [9]. Other applications might include using the Perceptron algorithm on problems where no kernel is available to exactly compute the weighted sums.

When Morris and Sinclair [28] solved the knapsack problem, they also generalized their result to a hypercube truncated by multiple hyperplanes (though the number of hyperplanes must be constant), rather than the single one that we consider in Section 4.2. Since the stationary distribution of their chain was assumed uniform, a natural question to ask is whether their results

can be generalized to non-uniform distributions, and if there are applications of this generalization to learning problems.

There is also the question of how to elegantly choose S and T for empirical use to balance time complexity and precision. While it is important to accurately estimate the weighted sums in order to properly simulate WM and Winnow, some imperfections in simulation can be handled since incorrect simulation decisions can be treated as noise, which Winnow and WM can tolerate. Ideally, the algorithms would intelligently choose S and T based on past performance, perhaps (for Winnow) utilizing the brute-force upper bound of $\alpha\theta$ on all weights in a brute-force execution (since no promotions can occur past that point). So $\forall \vec{p}, z(\vec{p}) \leq 1 + \lfloor \log_{\alpha} \theta \rfloor$. If this bound is exceeded during a run of Winnow, then one can increase S and T and run again.

Acknowledgments

The authors thank Mark Jerrum and Alistair Sinclair for their discussions, Jeff Jackson, Qingping Tao, and the COLT and JCSS reviewers for their helpful comments and Jeff Jackson for presenting an earlier version of this paper at COLT. This work was supported in part by NSF Grants CCR-0092761 and CCR-9877080 with matching funds from UNL-CCIS and a Layman grant, and was completed in part utilizing the Research Computing Facility of the University of Nebraska. Deepak Chawla performed this work at the University of Nebraska.

References

- [1] C. Blake, E. Keogh, C.J. Merz, UCI repository of machine learning databases, <http://www.ics.uci.edu/~mllearn/MLRepository.html> (2003).
- [2] A. Blum, P. Chalasani, J. Jackson, On learning embedded symmetric concepts, in: Proceedings of the Sixth Annual Workshop on Computational Learning Theory, ACM Press, New York, NY, 1993, pp. 337–346.
- [3] A. Blum, M. Furst, J. Jackson, M. Kearns, Y. Mansour, S. Rudich, Weakly learning DNF and characterizing statistical query learning using Fourier analysis, in: Proceedings of 26th ACM Symposium on Theory of Computing, 1994, pp. 253–262.
- [4] N.H. Bshouty, Simple learning algorithms using divide and conquer, *Comput. Complexity* 6 (2) (1997) 174–194.
- [5] N. H. Bshouty, J. Jackson, C. Tamon, More efficient PAC-learning of DNF with membership queries under the uniform distribution, *J. Comput. System Sci.*, to appear (early version in COLT 99).
- [6] N. Cesa-Bianchi, Y. Freund, D. Helmbold, D. Haussler, R. Schapire, M. Warmuth, How to use expert advice, *J. ACM* 44 (3) (1997) 427–485.
- [7] D. Chawla, L. Li, S.D. Scott, Efficiently approximating weighted sums with exponentially many terms, in: Proceedings of the 14th Annual Conference on Computational Learning Theory, 2001, pp. 82–98.
- [8] N. Cristianini, J. Shawe-Taylor, *An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods*, Cambridge, MA, Cambridge University Press, 2000.
- [9] T.G. Dietterich, R.H. Lathrop, T. Lozano-Perez, Solving the multiple-instance problem with axis-parallel rectangles, *Artificial Intelligence* 89 (1–2) (1997) 31–71.
- [10] M. Dyer, A. Frieze, R. Kannan, A. Kapoor, U. Vazirani, A mildly exponential time algorithm for approximating the number of solutions to a multidimensional knapsack problem, *Combin. Probab. Comput.* 2 (1993) 271–284.
- [11] Y. Freund, R.E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, *J. Comput. System Sci.* 55 (1) (1997) 119–139.
- [12] C. Gentile, M. K. Warmuth, Linear hinge loss and average margin, in: M.S. Kearns, S.A. Solla, D.A. Cohn (Eds.), *Advances in Neural Information Processing Systems*, Vol. 11, MIT Press, Cambridge, MA, 1998, pp. 225–231.

- [13] S.A. Goldman, S.K. Kwek, S.D. Scott, Agnostic learning of geometric patterns, *J. Comput. System Sci.* 6 (1) (2001) 123–151.
- [14] S.A. Goldman, S.D. Scott, Multiple-instance learning of real-valued geometric patterns, *Ann. Math. Artificial Intelligence* 39 (3) (2003) 259–290.
- [15] D.P. Helmbold, R.E. Schapire, Predicting nearly as well as the best pruning of a decision tree, *Mach. Learning* 27 (1) (1997) 51–68.
- [16] M. Jerrum, A. Sinclair, The Markov chain Monte Carlo method: an approach to approximate counting and integration, in: D. Hochbaum (Ed.), *Approximation Algorithms for NP-Hard Problems*, Boston, MA, PWS Pub., 1996, pp. 482–520 (Chapter 12).
- [17] M.R. Jerrum, L.G. Valiant, V.V. Vazirani, Random generation of combinatorial structures from a uniform distribution, *Theoret. Comput. Sci.* 43 (1986) 169–188.
- [18] R. Khairon, D. Roth, R. Servedio, Efficiency versus convergence of Boolean kernels for online learning algorithms, in: T.G. Dietterich, S. Becker, Z. Ghahramani (Eds.), *Advances in Neural Information Processing Systems*, Vol. 14, 2001, MIT Press, Cambridge, MA, pp. 423–430.
- [19] S. Kirkpatrick, C.D. Gelatt, M.P. Vecchi, Optimization by simulated annealing, *Science* 220 (1983) 671–680.
- [20] J. Kivinen, M.K. Warmuth, P. Auer, The perceptron algorithm vs. Winnow: linear vs. logarithmic mistake bounds when few input variables are relevant, *Artificial Intelligence* 97 (1–2) (1997) 325–343.
- [21] N. Littlestone, Learning quickly when irrelevant attributes abound: a new linear-threshold algorithm, *Mach. Learning* 2 (1988) 285–318.
- [22] N. Littlestone, From on-line to batch learning, in: *Proceedings of the Second Annual Workshop on Computational Learning Theory*, Morgan Kaufmann, Los Altos, CA, 1989, pp. 269–284.
- [23] N. Littlestone, Redundant noisy attributes, attribute errors, and linear threshold learning using Winnow, in: *Proceedings of the Fourth Annual Workshop on Computational Learning Theory*, Morgan Kaufmann, San Mateo, CA, 1991, pp. 147–156.
- [24] N. Littlestone, M.K. Warmuth, The weighted majority algorithm, *Inform. and Comput.* 108 (2) (1994) 212–261.
- [25] W. Maass, M.K. Warmuth, Efficient learning with virtual threshold gates, *Inform. Comput.* 141 (1) (1998) 66–83.
- [26] D.D. Margineantu, T.G. Dietterich, Pruning adaptive boosting, in: *Proceedings of the 14th International Conference on Machine Learning*, Morgan Kaufmann, Los Altos, CA, 1997, pp. 211–218.
- [27] N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, E. Teller, Equation of state calculation by fast computing machines, *J. Chem. Phys.* 21 (1953) 1087–1092.
- [28] B. Morris, A. Sinclair, Random walks on truncated cubes and sampling 0–1 knapsack solutions, *SIAM J. Comput.*, to appear (early version in FOCS 99).
- [29] K.-R. Muller, S. Mika, G. Ratsch, K. Tsuda, B. Schölkopf, An introduction to kernel-based learning methods, *IEEE Trans. Neural Networks* 12 (2) (2001) 181–201.
- [30] F. Pereira, Y. Singer, An efficient extension to mixture techniques for prediction and decision trees, *Mach. Learning* 36 (3) (1999) 183–199.
- [31] J.R. Quinlan, *C4.5: Programs for Machine Learning*, Morgan Kaufmann, Los Altos, CA, 1993.
- [32] J.R. Quinlan, Bagging, boosting, and C4.5, in: *Proceedings of the 13th National Conference on Artificial Intelligence*, 1996, pp. 725–730.
- [33] F. Rosenblatt, The perceptron: a probabilistic model for information storage and organization in the brain, *Psych. Rev.* 65 (1958) 386–407 (Reprinted in *Neurocomputing*, MIT Press, Cambridge, MA, 1988).
- [34] R.E. Schapire, Y. Freund, P. Bartlett, W.S. Lee, Boosting the margin: a new explanation for the effectiveness of voting methods, *Ann. Statist.* 26 (5) (1998) 1651–1686.
- [35] R.E. Schapire, Y. Singer, Improved boosting algorithms using confidence-rated predictions, *Mach. Learning* 37 (3) (1999) 297–336.
- [36] B. Schölkopf, A.J. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, MIT Press, Cambridge, MA, 2002.
- [37] S.D. Scott, J. Zhang, J. Brown, On generalized multiple-instance learning, Technical Report UNL-CSE-2003-5, Dept. of Computer Science, University of Nebraska, 2003.
- [38] A. Sinclair, Improved bounds for mixing rates of Markov chains and multicommodity flow, *Combin. Probab. Comput.* 1 (1992) 351–370.

- [39] E. Takimoto, M. Warmuth, Predicting nearly as well as the best pruning of a planar decision graph, *Theoret. Comput. Sci.* 288 (2) (2002) 217–235.
- [40] C. Tamon, J. Xiang, On the boosting pruning problem, in: *Proceedings of the 11th European Conference on Machine Learning*, Springer, Berlin, 2000, pp. 404–412.
- [41] Q. Tao, S. Scott, An analysis of MCMC sampling methods for estimating weighted sums in Winnow, in: C.H. Dagli (Ed.), *Artificial Neural Networks in Engineering*, ASME Press, Fairfield, NJ, 2003, pp. 15–20.
- [42] Q. Tao, S.D. Scott, A faster algorithm for generalized multiple-instance learning, in: *Proceedings of the Seventeenth Annual FLAIRS Conference*, AAAI Press, Miami Beach, FL, 2004, to appear.